

NEY YORK UNIVERSITY

CAPSTONE THESIS

**Dynamic equivalency and peer grouping
using Factor Modeling and Machine
Learning**

Author:
Jielan ZHENG, Yashu ZHU,
Jiangguangyu XUE

Supervisor:
Mr. Chris ROMANO

*A thesis submitted in fulfillment of the requirements
for the degree of Master of Finance and Risk Engineering*

August 13, 2019

NEY YORK UNIVERSITY

Abstract

Finance and Risk Department
Tandon School Of Engineering

Master of Finance and Risk Engineering

Dynamic equivalency and peer grouping using Factor Modeling and Machine Learning

by Jielan ZHENG, Yashu ZHU, Jiangguangyu XUE

This capstone project focus on cluster analysis of financial assets. Clustering is used as important knowledge discovery tools in modern machine learning process. The cluster analysis serves as a method to find which assets are different from each other. Given the wide spread use of factor modeling and machine learning, a firm should be able to dynamically create peer grouping. While diversification is a key step for constructing portfolios but true diversification is not possible as investing in each company requires a lot of capital. we discover the groupings of stocks which we get data from current stock market and transform data into organized structure through machine learning technics. With several clustering models and conduction of Well dividing stocks into groups with “similar characteristics” can help in portfolio construction to ensure we choose a universe of stocks with sufficient diversification between them....

Acknowledgements

We would like to extend our sincere gratitude to Mr. Chris Romano and Professor Agnes Tourin, for their expert guidance and constructive feedback throughout the course of this project. We appreciate the patience with which Mr. Chris would tackle our doubts and the time he would spare from his busy schedule to meet during the course of his work-day to discuss. His inputs to this project are invaluable. We would also like to thank the Finance and Risk Engineering department at NYU Tandon School of Engineering to give us the opportunity to work on this advanced topic for our capstone project, during the course of which we learned a lot. ...

Team August 13, 2019

Contents

Abstract	iii
Acknowledgements	v
1 Introduction	1
1.1 Introduction	1
2 Background	3
2.1 Feature Engineering	3
2.1.1 Feature Engineering in stock data	3
2.2 Principle Component Analysis	3
2.3 Clustering[5]	4
2.3.1 K-means and MiniBatchKmeans	4
2.3.2 Birch	5
2.3.3 Hierarchical clustering	5
2.3.4 SpectralClustering	6
2.3.5 DBSCAN	6
2.3.6 Mean Shift	6
2.3.7 Affinity Propagation	7
2.4 Distance Metrics	7
2.4.1 Euclidean Distance	9
2.4.2 Manhattan Distance	9
2.4.3 Chebychev Distance	9
2.4.4 Minkowski Distance	9
2.5 Clustering performance evaluation	9
2.5.1 Calinski-Harabasz score	9
2.5.2 Silhouette Coefficient score	10
2.5.3 Elbow method	10
2.5.4 Gap statistic	10
3 Data	13
3.1 Data Collection	13
3.2 Feature Engineering	13
3.2.1 Technical Features	13
3.2.2 Fundamental Features	14
3.2.3 Feature Combination	14
3.3 Principal Component Analysis (PCA)	16
3.4 Outlier Detection	18
4 Clustering Model	19
4.1 KMeans	19
4.1.1 PCA Effect Evaluation	19
4.1.2 Elbow Method	20

4.1.3	Gap Statistic Method	20
4.1.4	Prediction Distribution	20
4.1.5	Outliers Assignment	20
4.2	MiniBatchKMeans	20
4.2.1	Modeling Evaluation Exploration	22
4.2.2	Prediction Distribution	22
4.2.3	Outliers Assignment	22
4.2.4	Model Conclusion Of KMeans And MiniBatchkmeans	22
4.3	Birch	23
4.3.1	Searching For Optimal Hyper Parameters	23
4.3.2	PCA Effect Evaluation	23
4.3.3	Prediction Distribution	23
4.3.4	Outliers Assignment	25
4.3.5	Model Conclusion	25
4.4	Agglomerative	25
4.4.1	Dataset before PCA	25
4.4.2	Dataset after PCA, with outliers	25
4.4.3	Dataset after PCA, without outliers	28
4.4.4	Outlier Assignment	28
4.4.5	Model Conclusion	28
4.5	Spectral	30
4.5.1	Dataset before PCA	30
4.5.2	Dataset after PCA, with outliers	32
4.5.3	Dataset after PCA, without outliers	32
4.5.4	Outlier Assignment	32
4.5.5	Model Conclusion	32
4.6	Mean Shift	35
4.6.1	Dataset before PCA with outliers	35
4.6.2	Dataset before PCA without outliers	36
4.6.3	Dataset after PCA with outliers and without outliers	36
4.6.4	Model Conclusion	36
4.7	DBSCAN	37
4.8	Affinity Propagation	38
4.8.1	Dataset before PCA with outliers	38
4.8.2	Dataset before PCA without outliers	38
Classify samples without outliers	38	
Assign outliers back to clusters	40	
4.8.3	Dataset after PCA with outliers	40
4.8.4	Dataset after PCA with outliers	40
Classify samples without outliers	40	
Assign outliers back to clusters	42	
4.8.5	Model Conclusion	42
5	Model Comparison	43
5.1	The distribution of the results	43
5.2	Silhouette Coefficient and CH score Comparisons	43
5.3	PCA transformation effects and outliers removal effects	44
5.4	Outliers Reassignment	44
6	Conclusion and Future Research	45

A	Code: Data Preprocessing	47
A.1	Code: Get Data	47
A.2	Code: Feature Engineering	49
A.3	Code: Principal Components Analysis and Outlier Processing	52
B	Code: Model Implementation	55
B.1	Model: KMeans	55
B.2	Model: MiniBatchKMeans	58
B.3	Model: Birch	61
B.4	Model: Agglomerative	63
B.5	Model: Spectral	68
B.6	Model: Mean Shift	71
B.7	Model: Affinity Propagation	75
B.8	Model: DBSCAN	80

List of Figures

1.1 Dataframe Construction	1
1.2 Model Construction	2
2.1 KMeans and MiniBatchKmeans	4
2.2 Birch	5
2.3 DBSCAN	6
2.4 Meanshift	7
2.5 Affinity Propagation Example	8
2.6 Elbow	10
3.1 Feature Engineering for Technical Data	14
3.2 Processed Technical Features	14
3.3 Correlations Heat map	16
3.4 Cumulative Variance of Components	17
3.5 Name of Components	18
4.1 PCA Effect Evaluation	19
4.2 Elbow Method Of KMeans	20
4.3 Gap Statistic KMeans	21
4.4 Prediction Distribution Of KMeans	21
4.5 Three Evaluation Method for MiniBatchkmeans	22
4.6 Prediction Distribution Of MiniBatchkmeans	23
4.7 PCA Effect Evaluation Of Birch	24
4.8 Prediction Distribution Of Birch	24
4.9 Agglomerative Clustering scores(before PCA)	26
4.10 Dendrogram for Agglomerative Clustering Result (before PCA)	26
4.11 Histogram for Agglomerative Clustering Result (before PCA)	26
4.12 Agglomerative Clustering Scores (after PCA, with outliers)	27
4.13 Dendrogram for Agglomerative Clustering Result (after PCA, with outliers)	27
4.14 Histogram for Agglomerative Clustering Result (after PCA, with outliers)	27
4.15 Agglomerative Clustering scores (after PCA, without outliers)	28
4.16 Dendrogram for Agglomerative Clustering n=6, 13 (after PCA, without outliers)	29
4.17 Histogram for Agglomerative Clustering n=6, 13 (after PCA, without outliers)	29
4.18 Histogram for Agglomerative Clustering n=6, 13 (Placing back outliers)	30
4.19 Spectral Clustering Model scores (before PCA)	31
4.20 Histogram for Spectral Clustering n=4, 9 (before PCA)	31
4.21 Spectral Clustering Model scores (after PCA, with outliers)	32
4.22 Histogram for Spectral Clustering n=4, 8, 12 (after PCA, with outliers)	33
4.23 Spectral Clustering Model scores (after PCA, without outliers)	33

4.24 Histogram for Spectral Clustering (after PCA, without outliers)	34
4.25 Histogram for Spectral Clustering (Placing back outliers)	34
4.26 Mean Shift Algorithm scores(before PCA, with outliers)	35
4.27 Histogram for Mean Shift Clustering(before PCA, with outliers)	35
4.28 Mean Shift Algorithm scores(before PCA, without outliers)	36
4.29 Histogram for Mean Shift Clustering Result(before PCA, without outliers)	36
4.30 Mean Shift Clustering Result (After PCA, with outliers)	37
4.31 Mean Shift Clustering Result (After PCA, without outliers)	37
4.32 DBSCAN Clustering Result	37
4.33 Affinity Propagation Clustering Scores (Before PCA, with outliers) . .	38
4.34 Histogram for Affinity Propagation Clustering Result (Before PCA with outliers)	39
4.35 Affinity Propagation Clustering Scores (Before PCA, without outliers)	39
4.36 Histogram for Affinity Propagation Clustering Result (Before PCA without outliers)	40
4.37 Affinity Propagation Clustering Scores (After PCA, with outliers) . .	41
4.38 Histogram for Affinity Propagation Clustering Result (After PCA without outliers)	41
4.39 Affinity Propagation Clustering Scores (After PCA, without outliers) .	41
4.40 Histogram for Affinity Propagation Clustering Result (After PCA without outliers)	42

List of Tables

3.1	Fundamental Features and Explanation	15
4.1	Agglomerative Clustering Result Comparison	28
4.2	Spectral Model Parameters and Explanation	30
4.3	Spectral Clustering Result Comparison	33
5.1	Model Silhouette Score Comparison	43
5.2	Model CH Score Comparison	44

Chapter 1

Introduction

1.1 Introduction

The goal of this project is to conduct peer grouping of SP500 components using machine learning-based algorithms. This capstone project does the following several things coded in Python: It first identify and acquire the technical and fundamental data acquisition of SP500 components; built python framework for clustering; then conduct feature engineering to improve dataset by selecting representatively technical indicators; thirdly utilize PCA to perform dimension reduction on fundamental and technical features, then remove outlier dataset and built the final database; main part including conducting several clustering models with train and test sample from the database, then apply flexible evaluation approaches to relative models and analysis the clustering performance under each model. The program has to run in order as described.

The whole python framework design is shown in two main parts: Data and Clustering Model application. The data part is shown in the flow chart Figure 1.1.

In this project, the grouping stocks pool is SP500 components. To get the best results from the models and algorithms, it first involves getting the most out of the data for your algorithms to work with and so the data. This is the process and practice of feature engineering solves. Its main purposes are 1). Extract and reduce the high noise contained in the raw input. 2). To generate features containing sufficient and useful information needed for stock clustering. 3). Combine both technical features and fundamental features universally used in stock models.

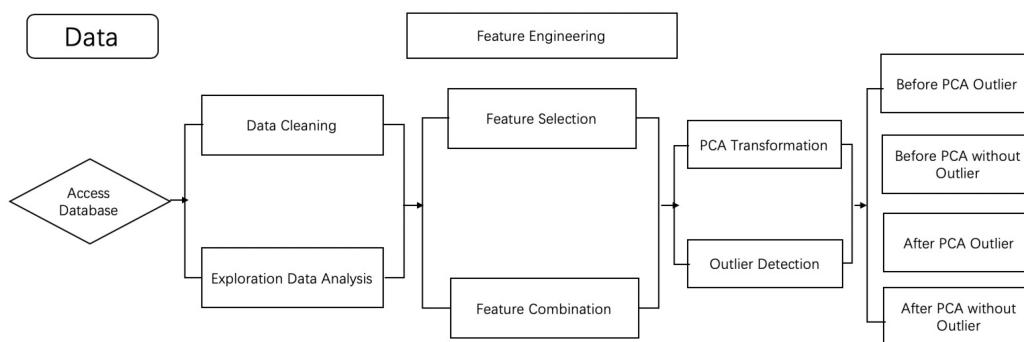


FIGURE 1.1: Dataframe Construction

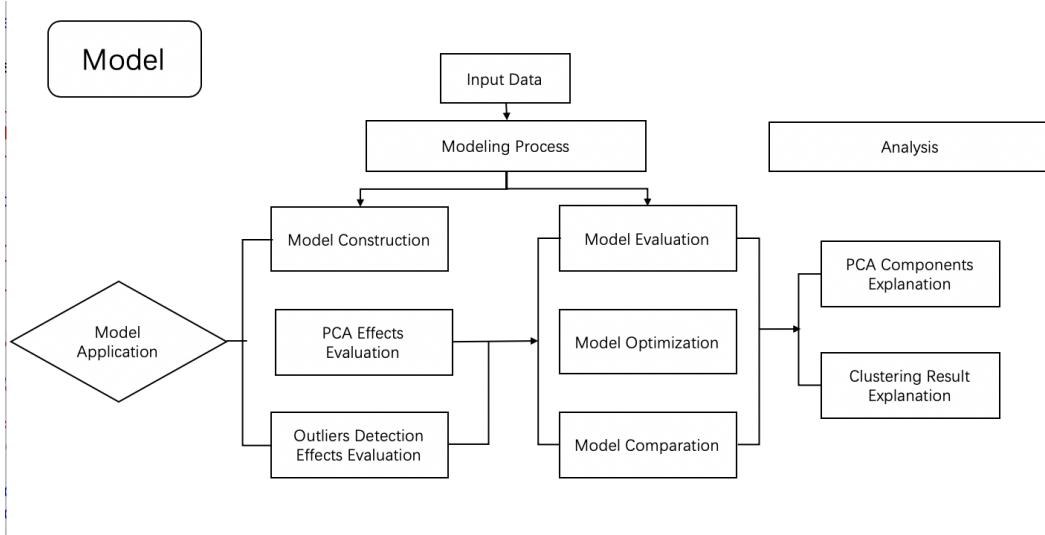


FIGURE 1.2: Model Construction

Moreover, we involved dimensionality reduction technique, Principal components analysis (PCA) and also remove outlier to achieve new database that 1). reduce dimensional feature-space 2). maintain as much information from the original dataset as possible in the reduced dataset. 3). exclude outlier noisy to the main sample.

The data part is shown in the flow chart Figure 1.2.

In the modeling process, we research and conduct several clustering models. To test the efficiency of PCA and then apply optimal datasets after PCA and excusalation of outliers as input to clustering models including K-means, MiniBatchKmeans, BIRCH, Agglomerative, Spectral, DBSCAN, Mean-Shift. To conduct properly valuation of performance under different clustering models, we work on several evaluation methods: 1). Calinski-Harabasz score 2). Silhouette Coefficient 3). Elbow Method 4). Gap Statistic Method.

Additionally, we manually assign the outliers into the “most similar” clustering group based on the distance calculation. Then compare model performances with full-dataset and dataset with outliers to test the distance metrics in grouping outliers and optimize final model results. Last, review and analysis all the model performances, draw the explanation of both PCA components analysis and peer grouping results.

We proceed as follows. In Chapter 2: Background, we provide some backgrounds in feature engineering, peer grouping model and machine learning techniques we used. Chapter 3: Data, we describe how we retrieve, parse and built data frame input. In Chapter 4: Clustering Model, we describe our modeling process, how we train, build and optimize our models, and relative evaluation implement included. In Chapter 5: Analysis we provide details of model performance and give an explanation of grouping results, and PCA factors analysis. In Chapter 6: Conclusion, we reach the conclusion based on all the analysis and reach and leave confusion and assumptions for further research. Last appendix where source codes will be provided.

Chapter 2

Background

2.1 Feature Engineering

Feature engineering is a term of art for data science and machine learning which refers to pre-processing and transforming raw data into a form which is more easily used by machine learning algorithms. Feature engineering is fundamental to the application of machine learning, and is both difficult and expensive. The need for manual feature engineering can be obviated by automated feature learning. This topic covers techniques for extracting and transforming features—the numeric representations of raw data—into formats for machine-learning models.[1]

Much like industrial processing can extract pure gold from trace elements within raw ore, feature engineering can extract valuable "alpha" from very noisy raw data.

2.1.1 Feature Engineering in stock data

Much research has been done in this field but it has been difficult to find the right features, which will help predict and classify with high accuracy. When you try to analyze stock classification using historical stock data, it has many attributes (feature set) containing useful as well as redundant features. Thus there is a need to remove the unwanted stock information.[2] In this research we first try to find the best features from the available company data as well as data about similar companies and stock indexes. The best-extracted features are then used to classify stock using different machine learning algorithms. [3]

2.2 Principle Component Analysis

Principle Component Analysis (PCA) is used for reducing the number of variables comprising a dataset while retaining the variability in the data, or identifying hidden patterns in the data, and classifying them according to how much of the information, stored in the data, they account for[4]. We use PCA for our stock market data for both purposes.

We take stock price panel data as score matrix A , and coefficient vector l that generates linear combination on A which yields principle components Y of smaller dimension than A . Principle components Y are independent and each coefficient vector l_i is required to maximize the variance of its corresponding principle component Y_i , as shown:

$$\max Var(Y_i) = l_i^T C l_i \quad (2.1)$$

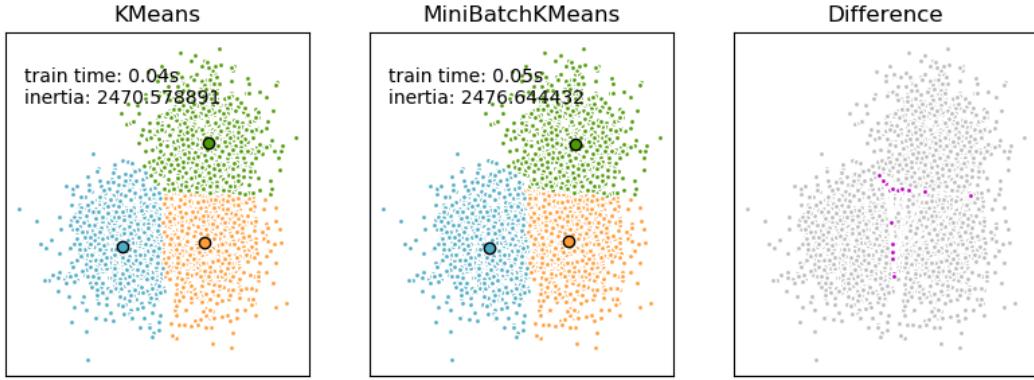


FIGURE 2.1: KMeans and MiniBatchKmeans

which has Lagrangian form:

$$\begin{aligned} s.t \quad & \|l_i^T\| = 1, \forall i \\ & l_i^T l_j = 0, \forall i \neq j \end{aligned}$$

Maximizing L by taking partial derivatives to 0, we obtain $C l_i = \lambda_i l_i$, where C is covariance matrix of A and λ_i, l_i are corresponding eigenvalues and eigenvectors. Therefore, principle components are constructed from eigenvectors of score matrix with score matrix itself: $Y_i = l_i^T A$.

2.3 Clustering[5]

2.3.1 K-means and MiniBatchKmeans

K-means is often referred to as Lloyd's algorithm. In basic terms, the algorithm has three steps. The first step chooses the initial centroids, with the most basic method being to choose k samples from the dataset X . After initialization, K-means consists of looping between the two other steps. The first step assigns each sample to its nearest centroid. The second step creates new centroids by taking the mean value of all of the samples assigned to each previous centroid. The difference between the old and the new centroids are computed and the algorithm repeats these last two steps until this value is less than a threshold. In other words, it repeats until the centroids do not move significantly.

The MiniBatchKMeans a variant of the KMeans algorithm which uses mini-batches to reduce the computation time, while still attempting to optimise the same objective function. Mini-batches are subsets of the input data, randomly sampled in each training iteration. These mini-batches drastically reduce the amount of computation required to converge to a local solution. In contrast to other algorithms that reduce the convergence time of k-means, mini-batch k-means produces results that are generally only slightly worse than the standard algorithm.

The algorithm iterates between two major steps, similar to vanilla k-means. In the first step, b samples are drawn randomly from the dataset, to form a mini-batch. These are then assigned to the nearest centroid. In the second step, the centroids are updated. In contrast to k-means, this is done on a per-sample basis. For each sample in the mini-batch, the assigned centroid is updated by taking the streaming average

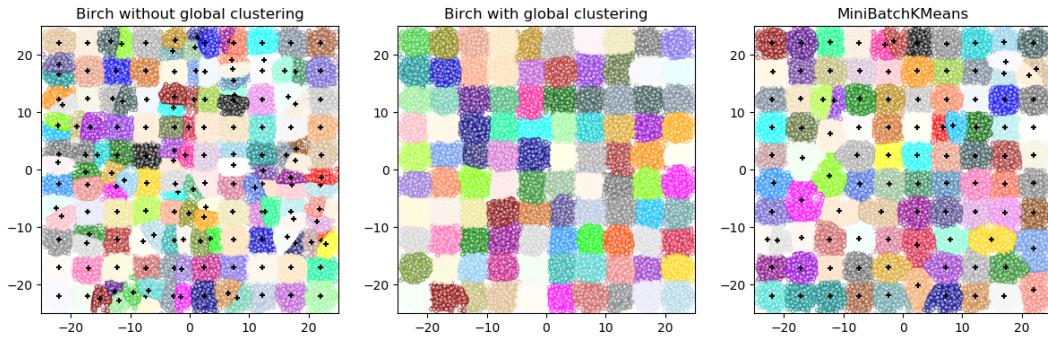


FIGURE 2.2: Birch

of the sample and all previous samples assigned to that centroid. This has the effect of decreasing the rate of change for a centroid over time. These steps are performed until convergence or a predetermined number of iterations is reached.

2.3.2 Birch

The Birch builds a tree called the Characteristic Feature Tree (CFT) for the given data. The data is essentially lossy compressed to a set of Characteristic Feature nodes (CF Nodes). The CF Nodes have a number of subclusters called Characteristic Feature subclusters (CF Subclusters) and these CF Subclusters located in the non-terminal CF Nodes can have CF Nodes as children. A new sample is inserted into the root of the CF Tree which is a CF Node. It is then merged with the subcluster of the root, that has the smallest radius after merging, constrained by the threshold and branching factor conditions. If the subcluster has any child node, then this is done repeatedly till it reaches a leaf. After finding the nearest subcluster in the leaf, the properties of this subcluster and the parent subclusters are recursively updated. If this split node has a parent subcluster and there is room for a new subcluster, then the parent is split into two. If there is no room, then this node is again split into two and the process is continued recursively, till it reaches the root. Birch does not scale very well to high dimensional data. As a rule of thumb if the number of feature is greater than twenty, it is generally better to use MiniBatchKMeans.

2.3.3 Hierarchical clustering

Hierarchical clustering is a general family of clustering algorithms that build nested clusters by merging or splitting them successively. This hierarchy of clusters is represented as a tree (or dendrogram). The root of the tree is the unique cluster that gathers all the samples, the leaves being the clusters with only one sample.

The AgglomerativeClustering object performs a hierarchical clustering using a bottom up approach: each observation starts in its own cluster, and clusters are successively merged together.

AgglomerativeClustering can also scale to large number of samples when it is used jointly with a connectivity matrix, but is computationally expensive when no connectivity constraints are added between samples: it considers at each step all the possible merges.

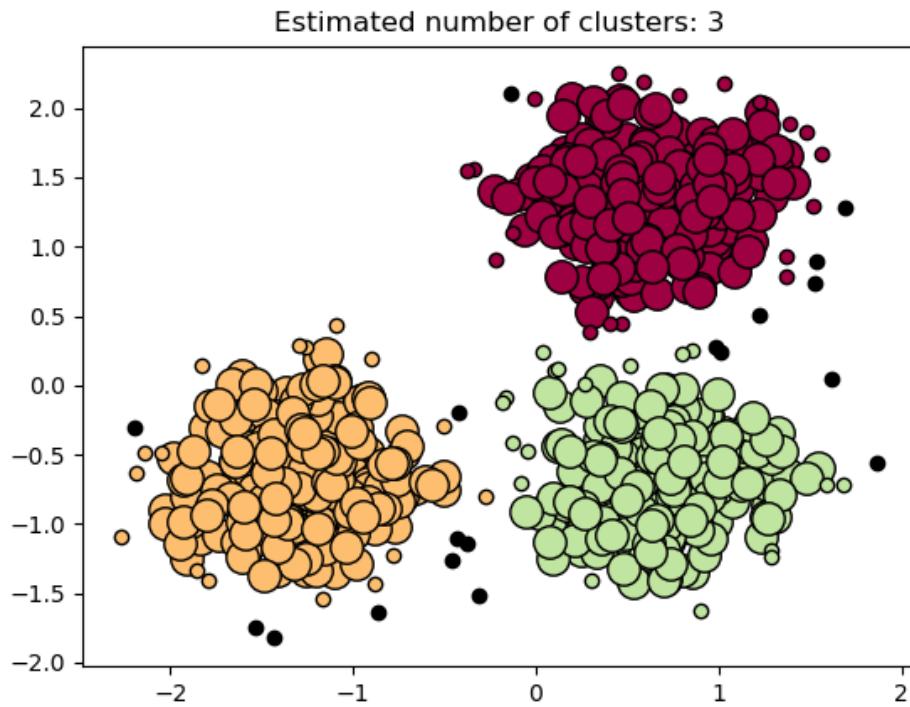


FIGURE 2.3: DBSCAN

2.3.4 SpectralClustering

SpectralClustering does a low-dimension embedding of the affinity matrix between samples, followed by a KMeans in the low dimensional space. It is especially efficient if the affinity matrix is sparse and the pyamgmodule is installed. SpectralClustering requires the number of clusters to be specified. It works well for a small number of clusters but is not advised when using many clusters.

2.3.5 DBSCAN

The DBSCAN algorithm views clusters as areas of high density separated by areas of low density. Due to this rather generic view, clusters found by DBSCAN can be any shape, as opposed to k-means which assumes that clusters are convex shaped. The central component to the DBSCAN is the concept of core samples, which are samples that are in areas of high density. A cluster is therefore a set of core samples, each close to each other (measured by some distance measure) and a set of non-core samples that are close to a core sample (but are not themselves core samples).

2.3.6 Mean Shift

MeanShift clustering aims to discover blobs in a smooth density of samples. It is a centroid based algorithm, which works by updating candidates for centroids to be the mean of the points within a given region. These candidates are then filtered in a post-processing stage to eliminate near-duplicates to form the final set of centroids.

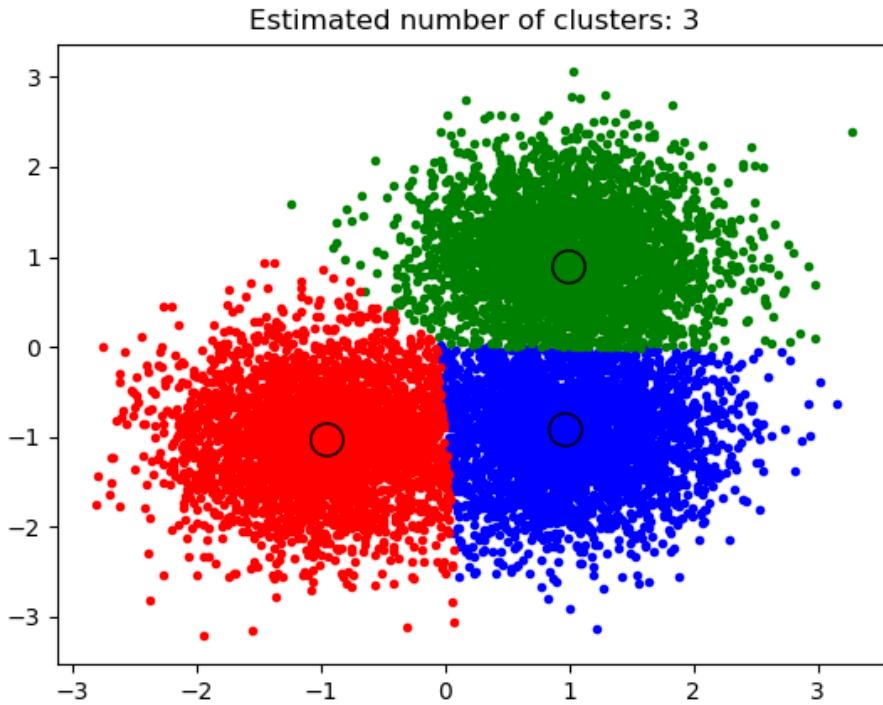


FIGURE 2.4: Meanshift

The algorithm is not highly scalable, as it requires multiple nearest neighbor searches during the execution of the algorithm. The algorithm is guaranteed to converge, however the algorithm will stop iterating when the change in centroids is small.

2.3.7 Affinity Propagation

Affinity Propagation creates clusters by sending messages between pairs of samples until convergence. A dataset is then described using a small number of exemplars, which are identified as those most representative of other samples. The messages sent between pairs represent the suitability for one sample to be the exemplar of the other, which is updated in response to the values from other pairs. This updating happens iteratively until convergence, at which point the final exemplars are chosen, and hence the final clustering is given.

2.4 Distance Metrics

Clustering is a technique to categorize the data into groups. Distance metrics plays a very important role in the clustering process. The more the similarity among the data in clusters, more the chances of particular data-items to belong to particular group. There are number of algorithms which are available for clustering.[6]

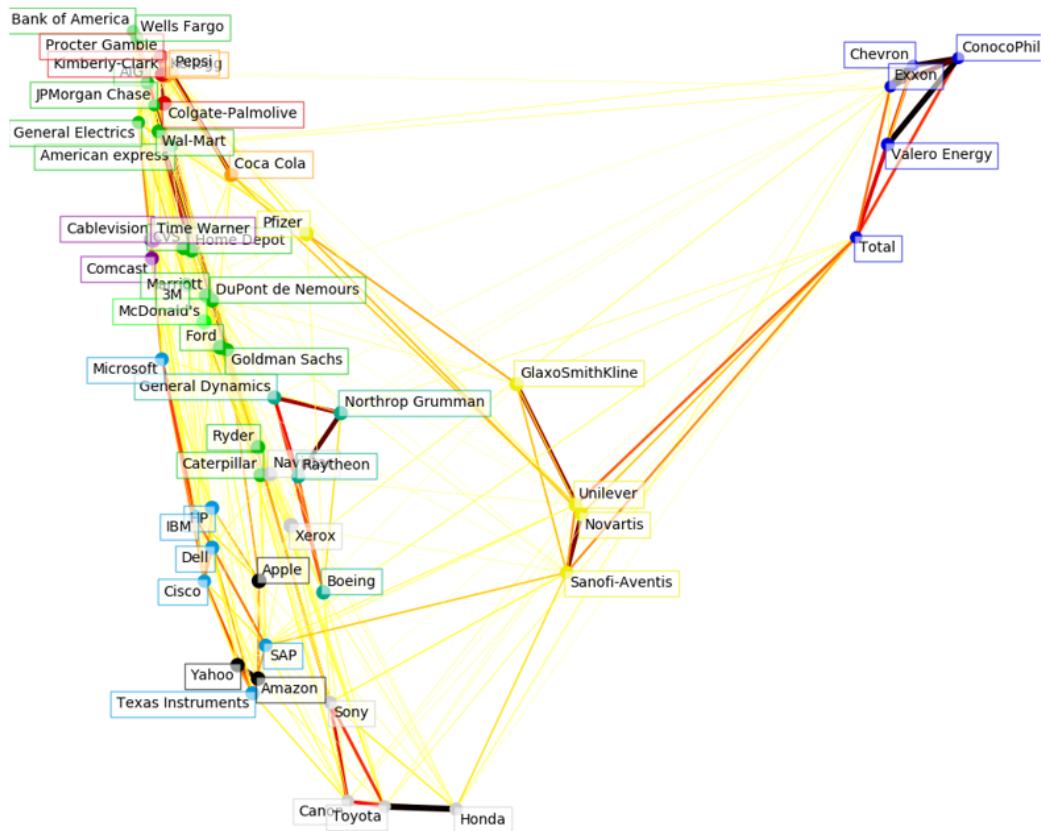


FIGURE 2.5: Affinity Propagation Example

2.4.1 Euclidean Distance

Euclidean distance computes the root of square difference between co-ordinates of pair of objects.

$$Dist_{XY} = \sqrt{\sum_{k=1}^d (X_{ik} - X_{jk})^2} \quad (2.2)$$

2.4.2 Manhattan Distance

Manhattan distance computes the absolute differences between coordinates of pair of objects

$$Dist_{XY} = |X_{ik} - X_{jk}| \quad (2.3)$$

2.4.3 Chebychev Distance

Chebychev Distance is also known as maximum value distance and is computed as the absolute magnitude of the differences between coordinate of a pair of objects

$$Dist_{XY} = \max |X_{ik} - X_{jk}| \quad (2.4)$$

2.4.4 Minkowski Distance

Minkowski Distance is the generalized metric distance

$$Dist_{XY} = \left[\sum_{k=1}^d |X_{ik} - X_{jk}|^{1/p} \right]^p \quad (2.5)$$

Note that when $p=2$, the distance becomes the Euclidean distance. When $p=1$ it becomes city block distance. Chebyshev distance is a variant of Minkowski distance where $p=\infty$ (taking a limit). This distance can be used for both ordinal and quantitative variables.

2.5 Clustering performance evaluation

Evaluating the performance of a clustering algorithm is not as trivial as counting the number of errors or the precision and recall of a supervised classification algorithm. In particular any evaluation metric should not take the absolute values of the cluster labels into account but rather if this clustering define separations of the data similar to some ground truth set of classes or satisfying some assumption such that members belong to the same class are more similar than members of different classes according to some similarity metric.[7]

2.5.1 Calinski-Harabasz score

If the ground truth labels are not known, the Calinski-Harabasz index can be used to evaluate the model, where a higher Calinski-Harabasz score relates to a model with better defined clusters. For k clusters, the Calinski-Harabasz score s is given as the ratio of the between-clusters dispersion mean and the within-cluster dispersion:

$$s(k) = \frac{Tr(B_k)}{Tr(W_k)} \times \frac{N - k}{k - 1}$$

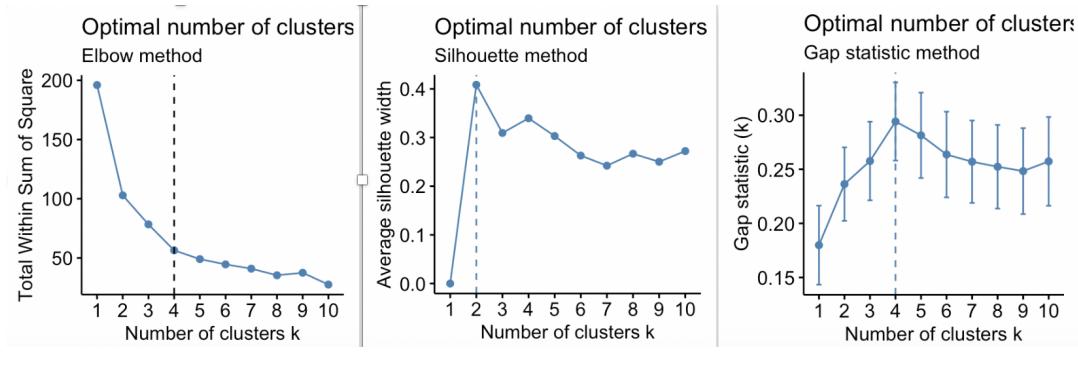


FIGURE 2.6: Elbow

where B_k is the between group dispersion matrix and W_k is the within-cluster dispersion matrix defined by:

$$W_k = \sum_{q=1}^k \sum_{x \in C_q} (x - c_q)(x - c_q)^T$$

$$B_k = \sum_q n_q (c_q - c)(c_q - c)^T$$

with N be the number of points in our data, C_q be the set of points in cluster q , c_q be the center of cluster q , c be the center of E , n_q be the number of points in cluster q .

2.5.2 Silhouette Coefficient score

If the ground truth labels are not known, evaluation must be performed using the model itself. Higher Silhouette Coefficient score relates to a model with better defined clusters. The Silhouette Coefficient is defined for each sample and is composed of two scores:

- a: The mean distance between a sample and all other points in the same class.
- b: The mean distance between a sample and all other points in the next nearest cluster. The Silhouette Coefficient s for a single sample is then given as:

$$s = \frac{b - a}{\max(a, b)}$$

2.5.3 Elbow method

The total WSS measures the compactness of the clustering and we want it to be as small as possible. The Elbow method looks at the total WSS as a function of the number of clusters: One should choose a number of clusters so that adding another cluster doesn't improve much better the total WSS. The location of a bend (knee) in the plot is generally considered as an indicator of the appropriate number of clusters.

2.5.4 Gap statistic

The gap statistic compares the total within intra-cluster variation for different values of k with their expected values under null reference distribution of the data.

The estimate of the optimal clusters will be value that maximize the gap statistic (i.e, that yields the largest gap statistic). This means that the clustering structure is far away from the random uniform distribution of points.[8]

Chapter 3

Data

We choose a dataset of stocks listed in Standard and Poor 500. The choice is made considering S&P 500 will provide our dataset proper diversity since the stocks span a variety of industries. According to our study, a well-diversified dataset is more suitable for clustering purpose [9]. The next step is to determine the appropriate features to use for implementation. Different financial datasets have been considered. One common choice for input is technical indicators of each stock. For instance, Wittman's analysis was based on three-year historical price data of 91 stocks [10]. Another team utilized the daily prices of S&P 500 stocks within a trading year as the input of their model [11]. Both the experiment generated reasonable results. Apart from technical features, fundamental data from financial statements are important data source for clustering analysis as well. Since all public companies are required to file financial reports periodically, fundamental indicators found in those reports should have relatively small noise and trustworthy data records [12]. In the experiment of Momeni et al., five fundamental ratios (ROA, ROE, Profit to sales ratio, EPS, and Operating profit margins) are selected for their K-means model [13]. Hence, both fundamental and technical features are considered for our clustering analysis.

3.1 Data Collection

All the quantitative data are collected from Yahoo Finance using a python library yahoofinancials. To get data from Yahoo Finance, we first get the stock tickers from Wikipedia. The technical part consists of daily prices (adjusted close, high, low) and volume of each S&P 500 listed stocks from Jun 2018 to Jun 2019. Technical data for S&P 500 index of the same time period is also retrieved for further feature engineering purpose. The fundamental data contains 51 key statistics data along with earnings per share, market cap, PE ratio, net income, and EBIT of each stocks. They will be further processed to generate the features we need for our model.

3.2 Feature Engineering

3.2.1 Technical Features

For technical data, there are three stocks (CTVA, DOW, ZTS) that have unmatching date variables with the others (the median of the group). We regard those as unqualified entries since it may cause errors when transforming them to pandas dataframe. Hence, those stock entries are dropped from the dataset. By researching relative case studies, considering both risk and return prospect, we choose the standard deviation of daily return and volume, the skewness and kurtosis of daily

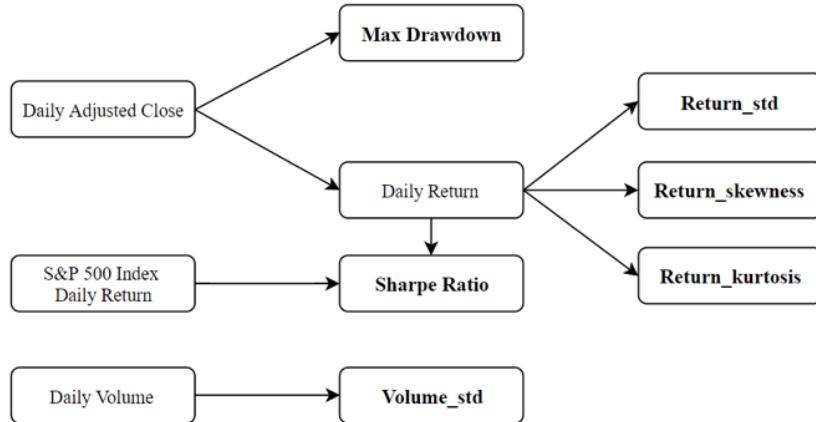


FIGURE 3.1: Feature Engineering for Technical Data

	Return_std	Return_skew	Return_kurtosis	Sharpe_Ratio	Volume_std	Maxdrawdown
ticker						
A	0.043125	0.572325	6.139398	0.029342	1.295404e+06	0.184895
AAL	0.059415	1.234017	8.127841	0.019145	3.349883e+06	0.386307
AAP	0.037719	-0.349844	9.141565	0.013477	5.052240e+05	0.305764
AAPL	0.048445	-0.566217	9.829872	0.015749	1.402302e+07	0.385159
ABBV	0.034534	0.708604	9.355761	-0.001167	2.907568e+06	0.206139

FIGURE 3.2: Processed Technical Features

return, Sharpe Ratio, and max drawdown as our technical features. Then, as Figure 3.1 suggested, they are calculated based on our processed technical data frame. These six factors are our technical input for clustering analysis (Figure 3.2).

3.2.2 Fundamental Features

For fundamental data, 57 indicators are collected originally. Regarding our existing dataset, here are two types of: (1) full of null entries (2) full of duplicated entries. Since no information is provided in those indicators, they are dropped out from our dataset. Then, to avoid comparison error in later steps, infinity entries are replaced by null value, which will be impute with mean or mode based on the pipeline setting later. After a series of processes, as Table 3.1 suggested, there are 28 fundamental features selected for model implementation.

3.2.3 Feature Combination

In this step, we combine the technical features and fundamental features into one single data frame for our clustering analysis. Due to the unmatching date (technical part) and the availability of financial statements (fundamental part) problem, two data frames are combined using inner join method. The integrated data frame is

TABLE 3.1: Fundamental Features and Explanation

Variables	Description
enterpriseToRevenue	Enterprise Value/Revenue, compares a company's enterprise value to its revenue
profitMargins	An indicator of a company's pricing strategies and how well it controls costs
enterpriseToEbitda	Enterprise Value/EBITDA, a popular valuation multiple used in the finance industry to measure the value of a company
52WeekChange	Movement of share prices in 52 weeks time frame
forwardEps	Earning per share calculated using Future earnings figures
sharesOutstanding	A company's stock currently held by all its shareholders
bookValue	The value of a security or asset as entered in a company's books
sharesShort	Amount of shares short selling
sharesPercentSharesOut	Short of Shares Outstanding
heldPercentInstitutions	Held by Institutions
netIncomeToCommon	Amount of net income to common shares
trailingEps	A company's earnings per share generated over a prior period
SandP52WeekChange	S&P500 52-Week Change
priceToBook	Price to Book ratio
heldPercentInsiders	Held by Insiders
shortRatio	Calculated amount of short interest that exists in a stock
floatShares	Shares of a public corporation that are available for trading in a stockmarket
Beta	Beta (3Y Monthly), a measure of volatility or systematic risk
enterpriseValue	A measure of a company's total value
earningsQuarterlyGrowth	An increase of a company's sales when compared to a previous quarter's revenue performance
pegRatio	A stock's price-to-earnings(P/E) ratio divided by the growth rate of its earnings for a specified time period
forwardPE	Calculated by as current stock price over the predicted next annual earnings period
shortPercentOfFloat	Short of Float
sharesShortPriorMonth	Shares Short (prior month)
eps	Portion of a company's profit allocated to each outstanding share of common stock
Market Cap	The value or capitalization the market puts on a company
PE	The ratio for valuing a company that measures its current share price relative to its per-share earnings
NI	The difference after factoring deductions and taxes into gross income
EBIT	Earnings before interest and tax

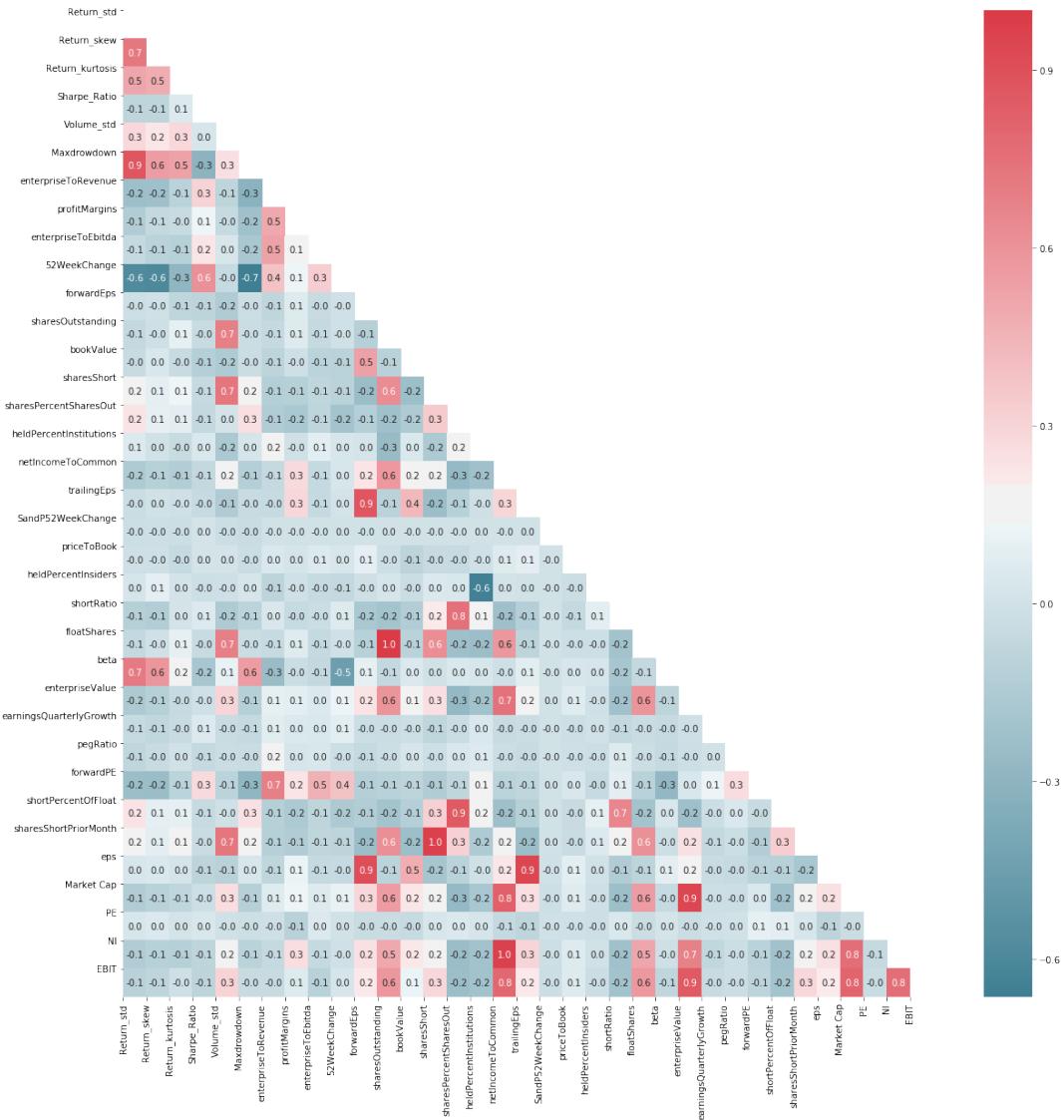


FIGURE 3.3: Correlations Heat map

then fed into our pipeline to impute missing values with corresponding mean and normalize all the entries. After such transformation, the data frame is qualified for dimensionality reduction.

3.3 Principal Component Analysis (PCA)

PCA is a common technique to reduce dimensions of features before training models. Since there are strong correlations among fundamental features, applying PCA here can help reduce the dimensions to a smaller number. And this can assist us in the future training of the clustering models because some of the models might work better with fewer dimensions.

As Figure 3.3 suggests, there are certain correlations among features we picked. This confirms our idea that there are strong correlations between fundamental features and PCA can be helpful here to reduce dimensions.

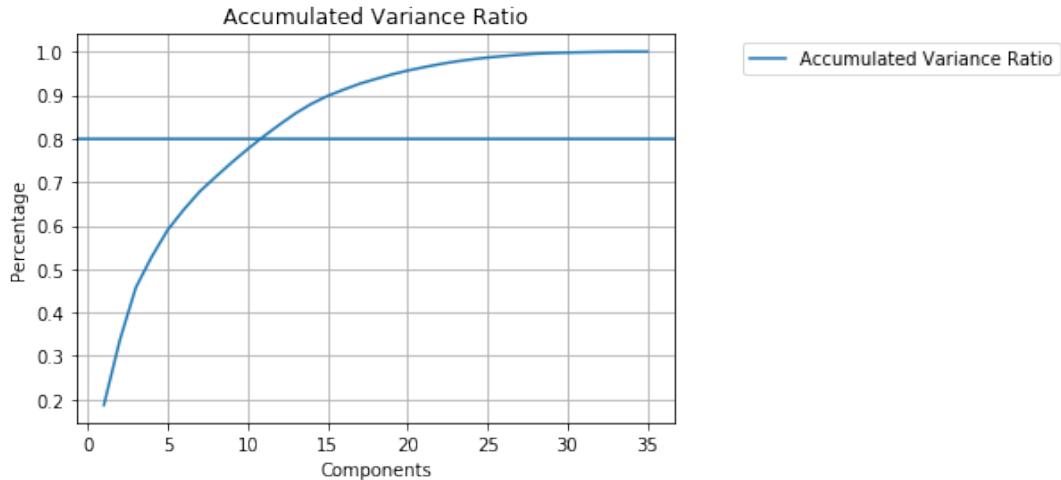


FIGURE 3.4: Cumulative Variance of Components

Then, we utilized PCA function in Scikit-Learn package to do the analysis and transformation of features.

To determine the number of dimensions we need to keep after transformation, the indicator of cumulative variance explained is focused in next stage. The components are sorted according to the cumulative variance explained and shown in the Figure 3.4.

As shown in Figure 3.4, if we keep top 11 components, we will have over 80 % variance explained. In other words, instead of using 35 dimensions to refer to 100 % information, we can just use 11 components to refer to 80 % of the information. That is efficient and satisfying.

However, we transform real world indicators (e.g. EPS, Return, PE) into 11 components with unknown meanings. To explore and get better understanding of the transformed components, we decide to dive into the loadings matrix to see the connections between components and original features.

To do that, for each component, in the loading matrix, we filter the features contribute over 0.3 or below -0.3. And we name the 11 components according to those original features that contribute more to the components. For example, according to loadings matrix,

$$\text{Component 2} = \text{Return_std} \times 0.32 + \text{Maxdrawdown} \times 0.33 + \text{Others}$$

There are only two features that have weights over 0.3 and these two features are related to risk. As a result, we can name component 2 as Risk Factor.

Similarly, we can name all the components according to loadings matrix (Figure 3.5).

10 of the 11 components can be given suitable name, while 8th component is hard to explain. The meaning is not explicit and we will just call it factor 8.

Generally speaking, after PCA transformation, we have factors related to Company Scale, Risk, EPS, Shorting ratio, Overall value to profitability, Holding Positions, Standard Deviation and so on.

Component	Name	Main constitutions	Component	Name	Main constitutions
1	Company Scale Factor	sharesOutstanding	0.32	6	Insiders Shares percents
		netIncomeToCommon	0.35		heldPercentInstitutions
		floatShares	0.32		heldPercentInstitutions
		enterpriseValue	0.34		Volume_std
		Market Cap	0.35		pegRatio
		NI	0.33		Sharpe_Ratio
2	Risk Factor	EBIT	0.34	8	Factor 8
		Return_std	0.32		priceToBook
		Maxdrawdown	0.33		PE
3	EPS Factor	forwardEps	0.36	9	Price to Profitability Ratio(The larger the worse)
		trailingEps	0.35		pegRatio
		eps	0.36		PE
4	Short Factor	sharesPercentSharesOut	0.44	10	Price to Book Ratio
		shortRatio	0.45		profitMargins
5	Overall value to Profitability (the larger the worse)	shortPercentOfFloat	0.43	11	Quarter Earnings Growth
		enterpriseToRevenue	0.41		priceToBook
		enterpriseToEbitda	0.35		Sharpe_Ratio
		forwardPE	0.35		earningsQuarterlyGrowth

FIGURE 3.5: Name of Components

3.4 Outlier Detection

Clustering method can be sensitive to outliers. A few outliers might cause clustering models to create unnecessary clusters. Consequently, dealing with outliers before fitting the models can be crucial.

We implemented outliers' detections and removals separately to features after PCA and features before PCA. We use 5-sigma rule to judge whether a sample is outlier: for a sample, if any of its feature is out of range of $[\mu - 5\sigma, \mu + 5\sigma]$, we classify it into outliers.

After that, we have four datasets: dataset after PCA with outliers, dataset after PCA without outliers, dataset before PCA with outliers and dataset before PCA without outliers.

These four datasets are used in different clusters. We compared the results to see whether PCA transformations and outliers detections improve models performance.

Chapter 4

Clustering Model

4.1 KMeans

The KMeans algorithm clusters data by trying to separate samples in n groups of equal variance, minimizing a criterion known as the inertia or within-cluster sum-of-squares. This algorithm requires the number of clusters to be specified. It scales well to large number of samples and has been used across a large range of application areas in many different fields.

And there are two main parameters we should assign before we fit the model. Set algorithm = full, which represents the classical EM-style algorithm and supports sparse data. The n_clusters presents the clustering number.

4.1.1 PCA Effect Evaluation

Test if PCA improve data structure in clustering model based on silhouette score(SH score) as evaluation method. Use dataset before PCA and after PCA, and modeling with optimal parameters. As presented in the graph below. The location of the maximum is considered as the appropriate number of clusters.

The evaluation result in the Figure 4.1 shows PCA improve the data structure and dataset after PCA has general higher SH score than that before PCA.

We continue the research dataset with after PCA. The optimal clustering number of After PCA is 3 and SH score is 0.207.

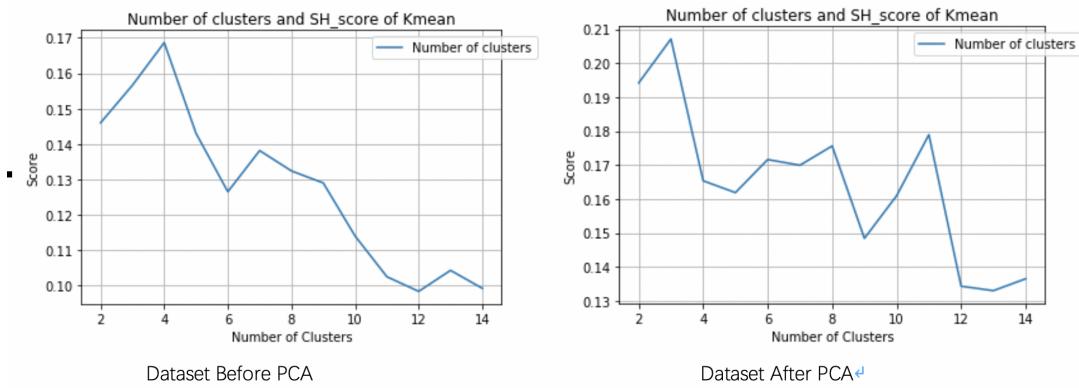


FIGURE 4.1: PCA Effect Evaluation

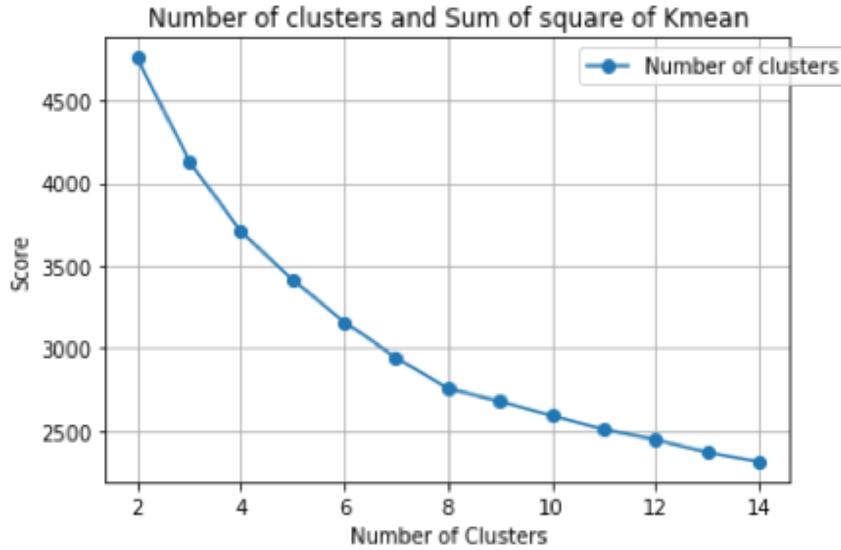


FIGURE 4.2: Elbow Method Of KMeans

4.1.2 Elbow Method

As mentioned in the background part, the location of a bend (knee) in the plot is generally considered as an indicator of the appropriate number of clusters. The curve in Figure 4.2 is not very explicit shows the optimal cluster number, and it only give a implicit range from 4 to 8.

4.1.3 Gap Statistic Method

The third evaluation method shows the results of the calculated gaps, the higher the value, the more optimal it is. We can see the relative optimal clusters number is 14 which is great different from the result above. But we could use the gap value of 3, 4, 5, 6, 8 as additional reference to evaluate other two metrics above. (Figure 4.3)

4.1.4 Prediction Distribution

Based on Figure 4.3 we see 8 and 11 as relative optimal clustering numbers. And the distribution graphs (Figure 4.4) show the prediction distribution of 8 and 11 clusters are both balanced.

4.1.5 Outliers Assignment

We remove outlier and find the relative clustering numbers above. Next we reassign those outlier based on the distance between each outlier and each cluster center. In the code part you can see the distance metric we use. After assignment, the SH score decrease as we expected and we do the assignment analyst in Chapter4.

4.2 MiniBatchKMeans

The MiniBatchKMeans is a variant of the KMeans algorithm which uses mini-batches to reduce the computation time, while still attempting to optimise the same

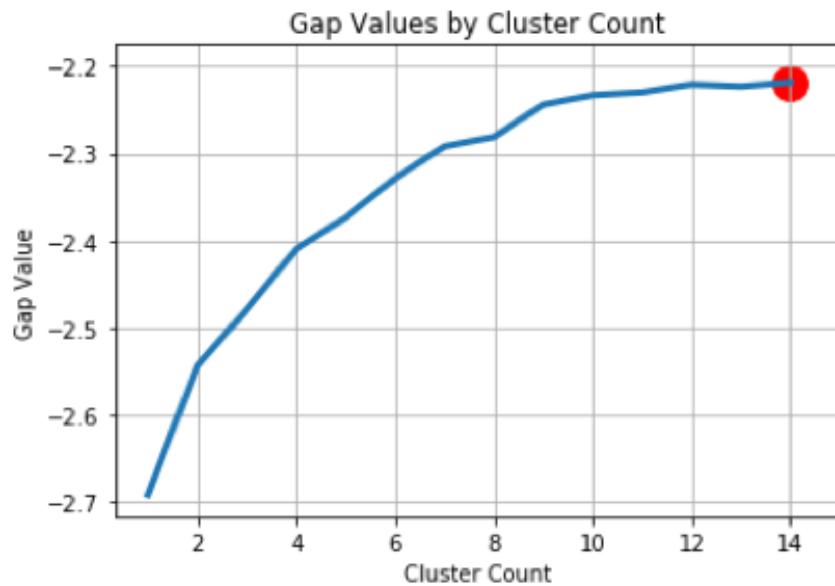


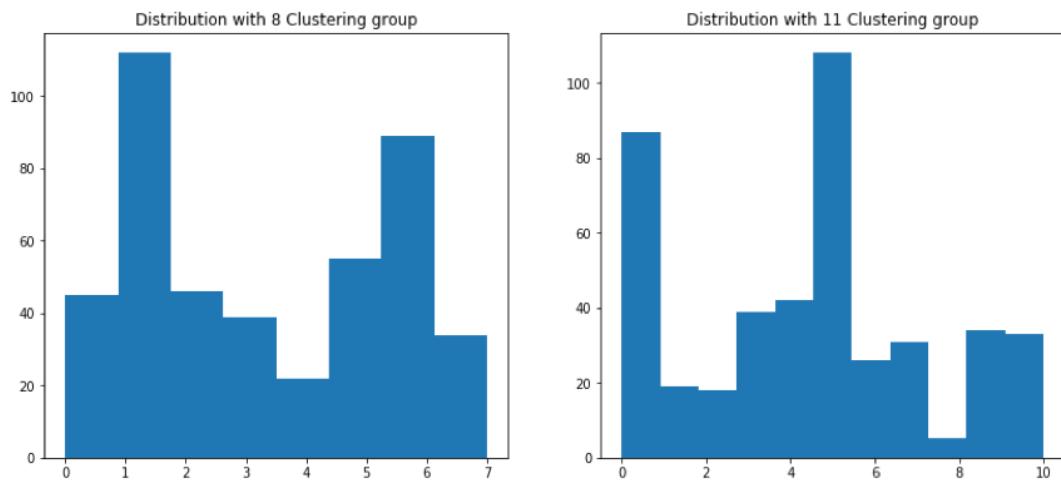
FIGURE 4.3: Gap Statistic KMeans

FIGURE 4.4: Prediction Distribution Of KMeans

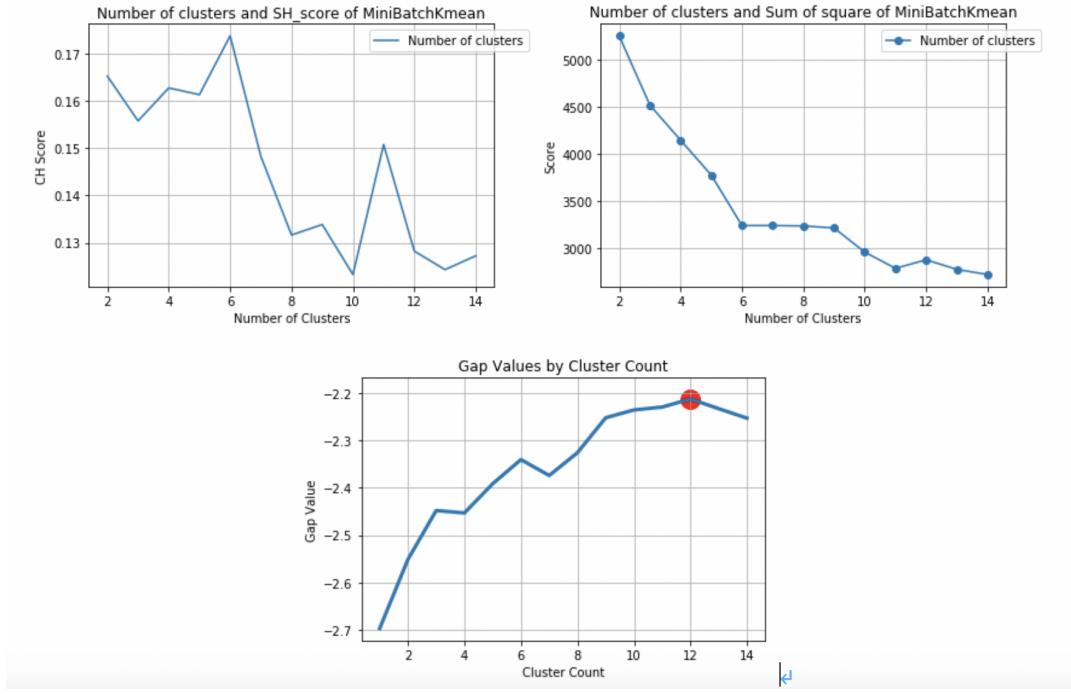


FIGURE 4.5: Three Evaluation Method for MiniBatchkmeans

objective function. Mini-batches are subsets of the input data, randomly sampled in each training iteration. So for this model we just briefly present all the results below.

4.2.1 Modeling Evaluation Exploration

The model performance based on three evaluation method (Figure 4.5).

4.2.2 Prediction Distribution

For MinibatchKmeans we have a more clear result. We could see 6 and 11 as good cluster numbers and we use that three to see the prediction result (Figure 4.6). We could have a relative balance result.

4.2.3 Outliers Assignment

And the last, reassign outliers back to groups based on the distance metrics the same as KMeans. Detailed Analysis presented in the see next Chapter.

4.2.4 Model Conclusion Of KMeans And MiniBatchkmeans

KMeans model works well on the dataset after PCA and without outlier, and present a relative balanced prediction distribution with 8 and 11 clusters.

While MinibatchKmeans, based on the same algorithm as Kmeans, optimising the same objective function, shows more clear result of performance evaluation. And it presents a relative balanced prediction distribution with 6 and 11 clusters.

And it worth to mention that comparing and MiniBatchKMeans the former always have better performance under all the evaluation scores. We could see that

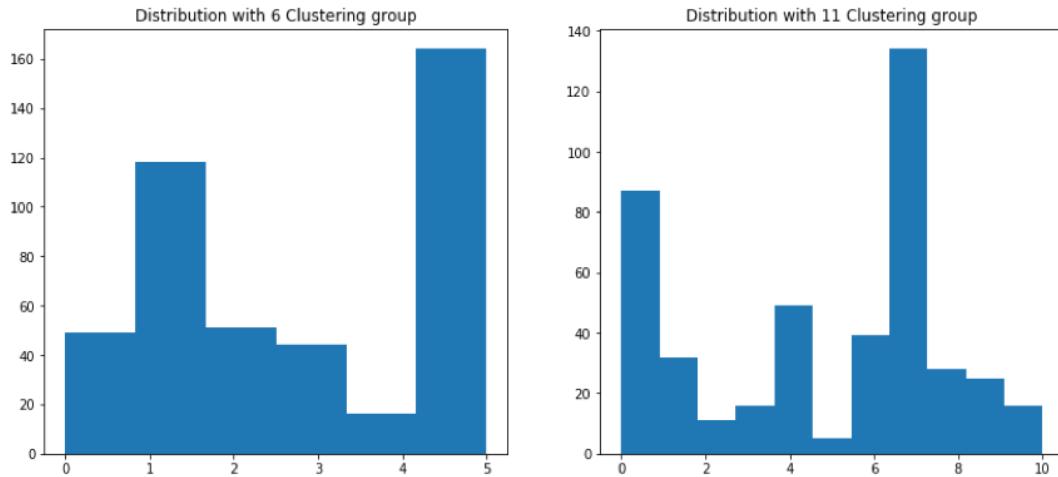


FIGURE 4.6: Prediction Distribution Of MiniBatchkmeans

MiniBatchKMeans converges faster than KMeans, but the quality of the results is reduced.

4.3 Birch

4.3.1 Searching For Optimal Hyper Parameters

The Birch algorithm has two parameters, the threshold and the branching factor. The branching factor limits the number of subclusters in a node and the threshold limits the distance between the entering sample and the existing subclusters.

Conduct Grid-search to find optimal parameters and then see the clustering result based on Calinski Carabaz score(CH score) performance evaluation. The Grid-search function code attached in appendix.

4.3.2 PCA Effect Evaluation

Based on CH score evaluation method we could also test if PCA improve data structure in Birch model. Figure 4.7 show optimal result of after PCA dataset and before PCA dataset.

In each case with same parameters, dataset after PCA has general higher CH scores than that before PCA. And Figure 4.7 only show the optimal parameter case of those two datasets. So we could conclude that PCA improve the data structure for Birch model.

4.3.3 Prediction Distribution

We continue the research dataset with after PCA. The relative optimal clustering number of after PCA is 4 and 8 with threshold = 1 and branchingfactor = 50. And the distribution graphs (Figure 4.8) show the prediction of 8 clusters has a more balanced distribution than that with 4 clusters.

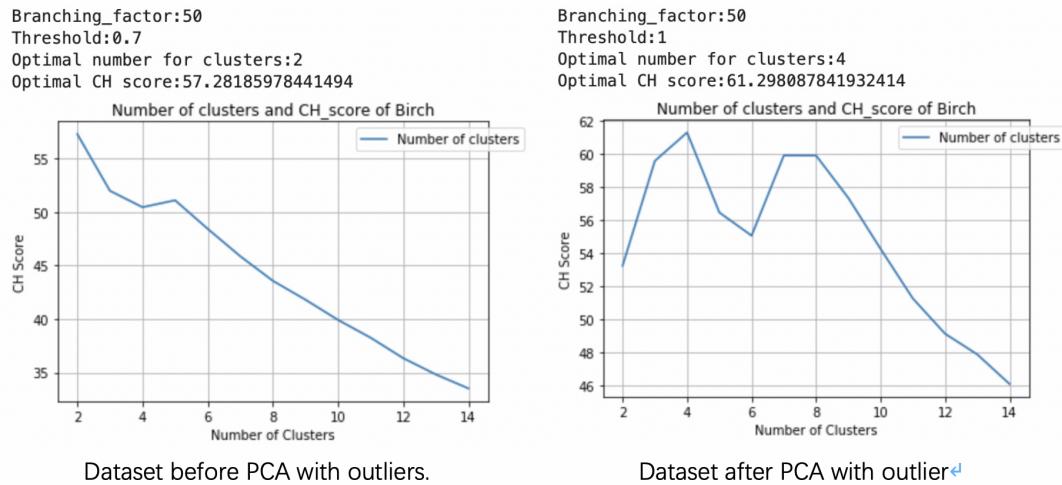


FIGURE 4.7: PCA Effect Evaluation Of Birch

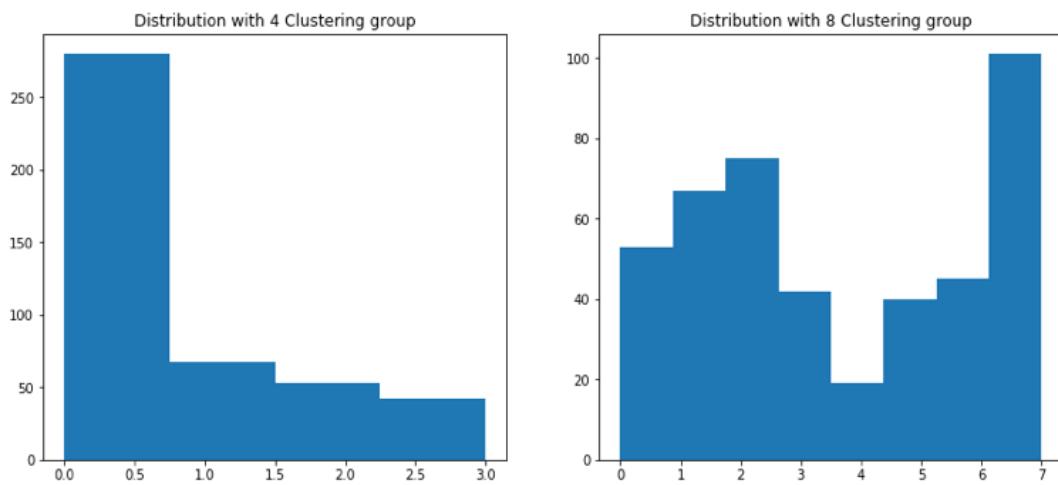


FIGURE 4.8: Prediction Distribution Of Birch

4.3.4 Outliers Assignment

We remove outlier and find the relative clustering numbers above. Next we reassign those outlier based on the distance between each outlier and each cluster center. After assignment, the CH score decrease, which is the same as KMeans and MinibatchKMeans.

4.3.5 Model Conclusion

Birch does not scale very well to high dimensional data. As a rule of thumb if features is greater than twenty, it is generally better to use MiniBatchKMeans. If the number of instances of data needs to be reduced, or if one wants a large number of subclusters either as a preprocessing step or otherwise, Birch is more useful than MiniBatchKMeans.

4.4 Agglomerative

For agglomerative clustering, there are three parameters could be changed: number of clusters, linkage type ('ward', 'complete', 'average', 'single'), and distance calculation methods ('euclidean', 'l1', 'l2', 'manhattan', 'cosine', 'precomputed'). We performed grid search for each type of dataset based on these three parameters with number of clusters changing from 2 to 20. For each cluster number, we find out the parameter combination that generates the highest Calinski-Harabasz (CH) score, since it is regarded as the primary evaluation method for agglomerative clustering model. A plot of CH score and Silhouette (SH) score versus number of clusters will be generate for each dataset. Moreover, SH scores are rescaled by multiply 100 for more intuitive comparison purpose.

4.4.1 Dataset before PCA

The result suggests that for dataset before PCA, the optimal combination for all cluster numbers happen to be linkage = 'ward' and affinity = 'euclidean'. As figure 4.9 suggested, the CH score for this dataset is descending as number of clusters increase. According to [14], descending or ascending or smooth-horizontal CH score doesn't provide any reference value – we cannot prefer one solution to others based on a line like this. Hence, we choose the optimal number of clusters through finding the peak for SH score.

As Figure 4.9 suggested, there are two abrupt elbows in SH score: cluster number equals to 3 and 12. Since we want an outcome with around 10 clusters, the optimal number of clusters for this dataset is therefore 12. According to Figure 4.10 and 4.11, it can be seen that the clusters are imbalanced.

4.4.2 Dataset after PCA, with outliers

With dimensionality reduction, the optimal configuration is still linkage = 'ward' and affinity = 'euclidean'. As Figure 4.12 suggested, the optimal number of clusters will be 9 for dataset after PCA. The imbalanced clusters problem still occurs (Figure 4.13, 4.14). One thing can be noticed is that the outcome after PCA is overall better than results before PCA: CH score increases around 40 % for all number of clusters.

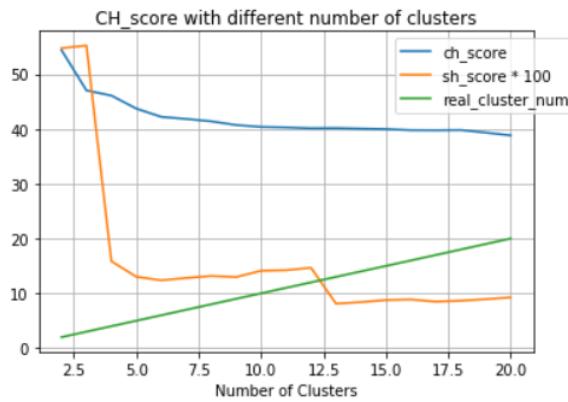


FIGURE 4.9: Agglomerative Clustering scores(before PCA)

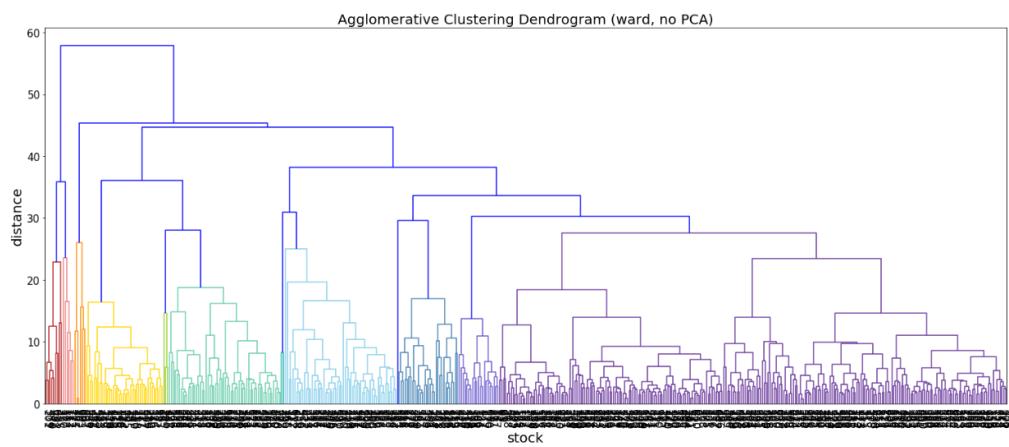


FIGURE 4.10: Dendrogram for Agglomerative Clustering Result (before PCA)

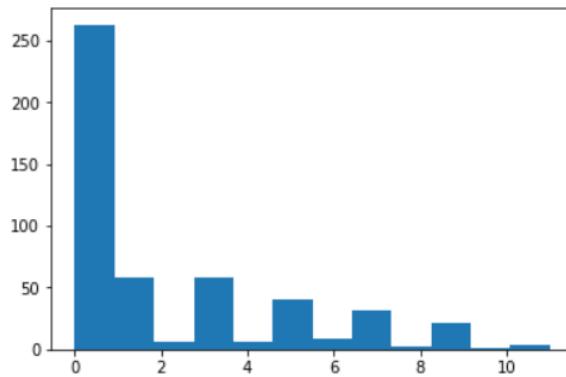


FIGURE 4.11: Histogram for Agglomerative Clustering Result (before PCA)

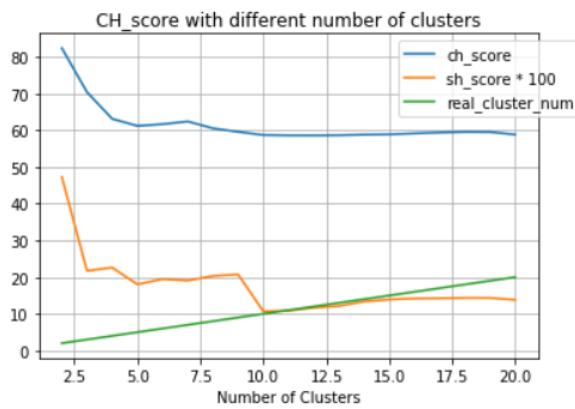


FIGURE 4.12: Agglomerative Clustering Scores (after PCA, with outliers)

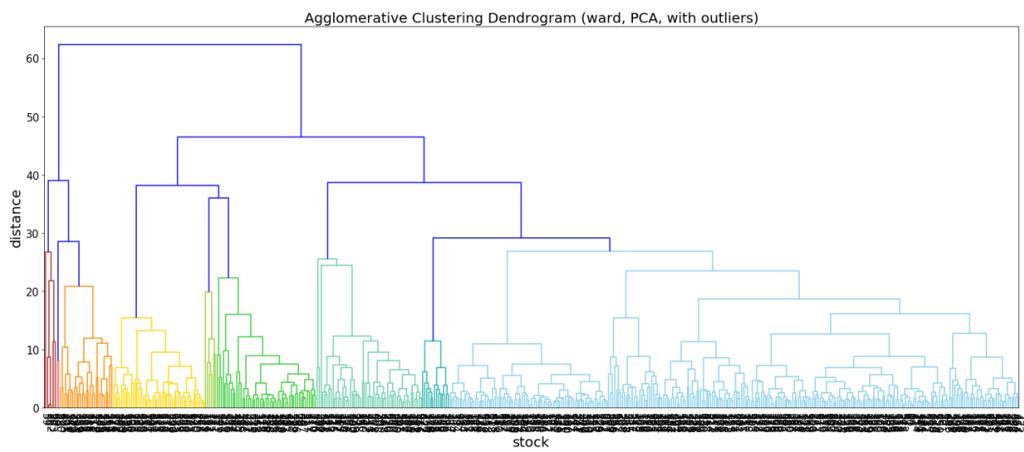


FIGURE 4.13: Dendrogram for Agglomerative Clustering Result (after PCA, with outliers)

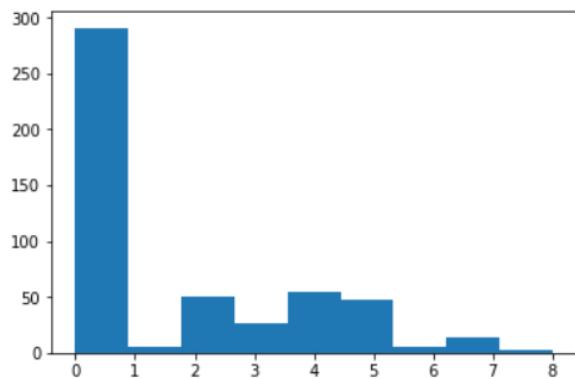


FIGURE 4.14: Histogram for Agglomerative Clustering Result (after PCA, with outliers)

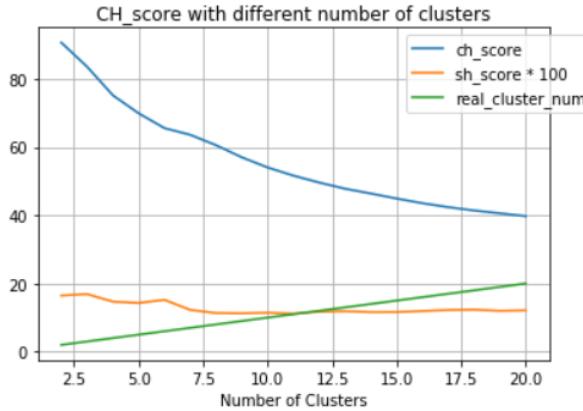


FIGURE 4.15: Agglomerative Clustering scores (after PCA, without outliers)

TABLE 4.1: Agglomerative Clustering Result Comparison

	# Clusters	CH score	SH score	Cluster_std
Before PCA (with outliers)	12	40.1227	0.1466	72.8629
After PCA (with outliers)	9	59.5859	0.2073	90.7277
After PCA (without outliers)	6	65.5629	0.1521	49.0252
After PCA (without outliers)	13	47.7946	0.1187	23.7171

4.4.3 Dataset after PCA, without outliers

After removing outliers, the combination of linkage = 'ward' and affinity = 'euclidean' still outperform the other combinations. While no abrupt elbow is noticed, there is a small peak for SH score when cluster number = 6 and 13 (Figure 4.15).

Both results generate relatively balanced clusters (Figure 4.16, 4.17).

4.4.4 Outlier Assignment

Then, outliers are labeled and added back to the whole dataset based on the clustering results of above two models. The clusters generated are still relatively balanced compared to the result before PCA and after PCA with outliers (Figure 4.18).

4.4.5 Model Conclusion

A more intuitive comparison can be seen from Table 4.1. We can conclude, for agglomerative clustering, using dataset after PCA without outliers will keep CH and SH score in acceptable range while generate relatively balanced clustering results. While it may be computationally expensive, Agglomerative clustering performs well with large number of samples and clusters. Moreover, this algorithm is also commonly used for unsupervised learning for its unique computation mechanism. In each successive iteration, it agglomerates the closest pair of samples by satisfying certain similarity criteria [15]. Because of this, outliers have a significant impact on the results. That is also the reason why after removing the outliers, the model generates more balanced clusters.

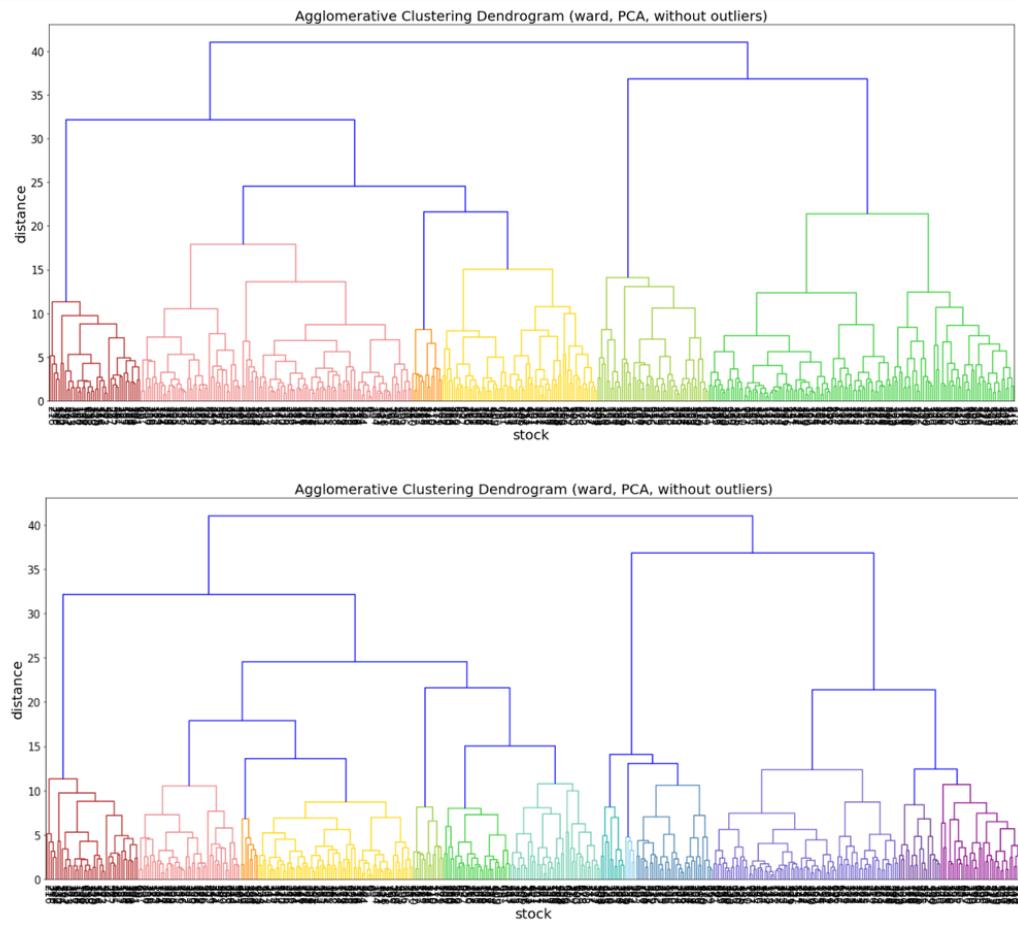


FIGURE 4.16: Dendrogram for Agglomerative Clustering n=6, 13 (after PCA, without outliers)

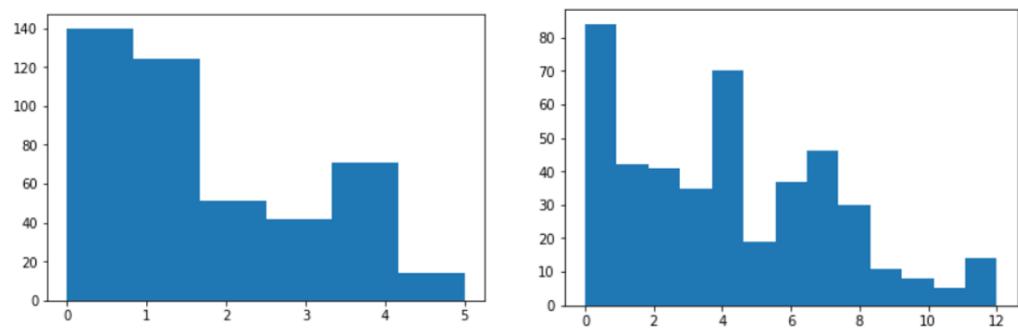


FIGURE 4.17: Histogram for Agglomerative Clustering n=6, 13 (after PCA, without outliers)

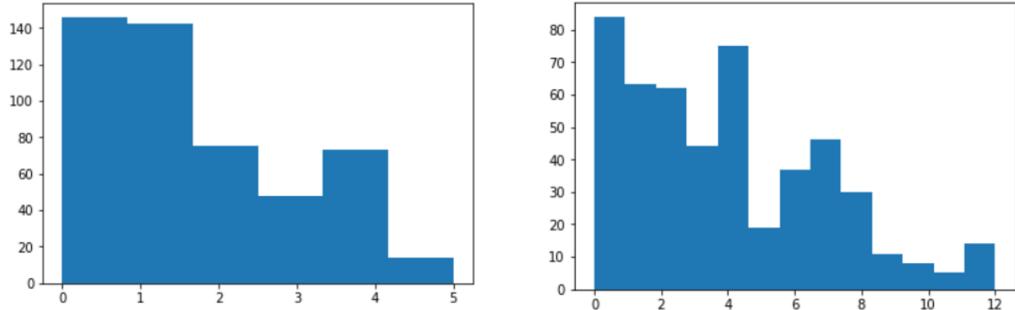


FIGURE 4.18: Histogram for Agglomerative Clustering n=6, 13 (Placing back outliers)

TABLE 4.2: Spectral Model Parameters and Explanation

Name	Description
n_clusters	Number of clusters
n_init	Number of time the k-means algorithm will be run with different centroid seeds. <i>np.arange(5,30,3)</i>
n_neighbors	Number of neighbors to use when constructing the affinity matrix using the nearest neighbors method. <i>np.arange(5,30,5)</i>
Gamma	Kernel coefficient for rbf, poly, sigmoid, laplacian and chi2 kernels. <i>np.arange(0.6,1.5,0.2)</i>
affinity	'nearest_neighbors', 'precomputed', 'rbf'

4.5 Spectral

For Spectral Clustering, there are five parameters that can be changed (Table 4.2). We tested the algorithm by ranging number of clusters from 2 to 20. As the customized grid search function for Agglomerative Clustering, we traversing all combinations to find out the one that generates the highest CH score for certain number of clusters.

4.5.1 Dataset before PCA

For dataset before PCA, there is a peak in both CH and SH score when number of clusters equals to 4 (Figure 4.19), $n_init = 5$, $\gamma = 0.6$, $affinity = \text{nearest_neighbors}$, $n_neighbors = 10$. This might because spectral clustering performs the best when given medium number of samples and small number of clusters [5]. However, the clustering result (Figure 4.20) suggests that the clusters generated are imbalanced.

To satisfy our clustering analysis purpose, more clusters is desired as we mentioned in the beginning of this thesis. We found another peak in SH score when number of clusters = 9, $n_init = 10$, $\gamma = 0.6$, $affinity = \text{nearest_neighbors}$, and $n_neighbors = 10$. The clusters are slightly more balanced than number of clusters equals to 4(Figure 4.20).

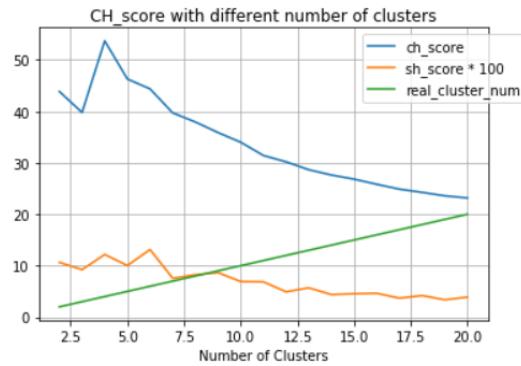


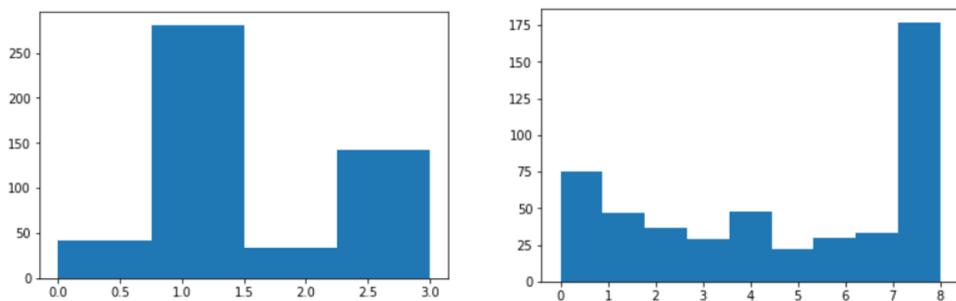
FIGURE 4.19: Spectral Clustering Model scores (before PCA)

FIGURE 4.20: Histogram for Spectral Clustering n=4, 9 (before PCA)

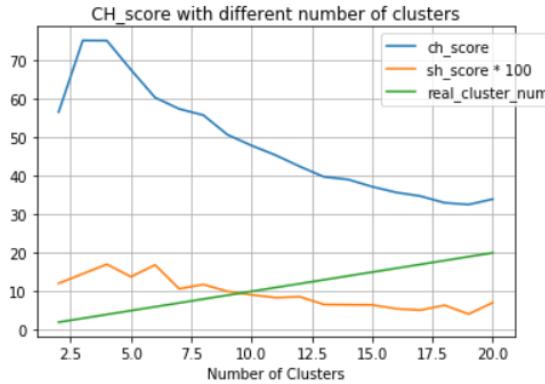


FIGURE 4.21: Spectral Clustering Model scores (after PCA, with outliers)

4.5.2 Dataset after PCA, with outliers

After PCA and before dealing with outliers, optimal combinations vary for different number of clusters. The overall CH scores increased compared with the outcome before PCA (Figure 4.21).

As above situation, a peak for both CH and SH score can be noticed when number of clusters = 4, n_init = 5, gamma = 0.6, affinity = nearest_neighbors, and number of neighbors = 10. In this case (Figure 4.22), the clusters are relatively imbalanced. Another two sets of configurations are also regarded as acceptable as they capture the peak in SH score. That is when number of clusters = 8, n_init = 20, gamma = 0.6, affinity = nearest_neighbors, number of neighbors = 15 and when number of clusters = 12, n_init = 10, gamma = 0.6, affinity = nearest_neighbors, number of neighbors = 15. As the number of clusters increased, the clustering outcome is getting more balanced.

4.5.3 Dataset after PCA, without outliers

After removing the outliers from dataset with PCA, the CH score increased significantly for small number of clusters (Figure 4.23). One small peak in both CH and SH score appears when number of clusters = 8, n_init = 5, gamma = 0.6, affinity = nearest_neighbors, and number of neighbors = 25.

The result is beautifully balanced (Figure 4.24).

4.5.4 Outlier Assignment

After labelling and adding back the outliers, the result is still balanced (Figure 4.25). Most of the outliers fall in the first three clusters. The overall shape of the clustering distribution doesn't change obviously.

4.5.5 Model Conclusion

A more intuitive comparison can be seen from Table 4.3. The dataset after PCA without outliers performs the best as well. We can also conclude that generally, Spectral Clustering algorithm works better for small number of clusters.

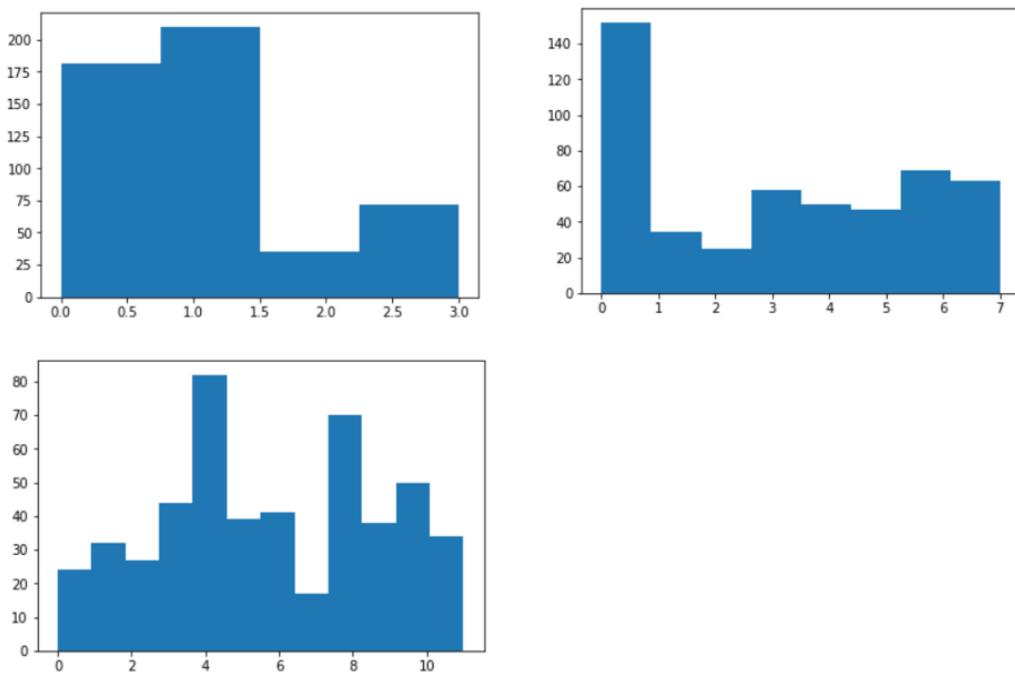


FIGURE 4.22: Histogram for Spectral Clustering n=4, 8, 12 (after PCA, with outliers)

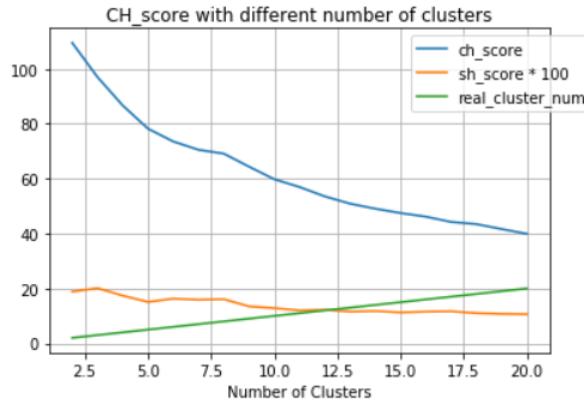


FIGURE 4.23: Spectral Clustering Model scores (after PCA, without outliers)

TABLE 4.3: Spectral Clustering Result Comparison

	# Clusters	CH score	SH score	Cluster_std
Before PCA	4	53.7176	0.1218	115.4369
Before PCA	9	35.8770	0.0864	48.2105
After PCA (with outliers)	4	75.0881	0.1703	84.2002
After PCA (with outliers)	8	55.7545	0.1179	39.0741
After PCA (with outliers)	12	42.4231	0.0860	18.6182
After PCA (without outliers)	8	69.1317	0.1609	24.9967

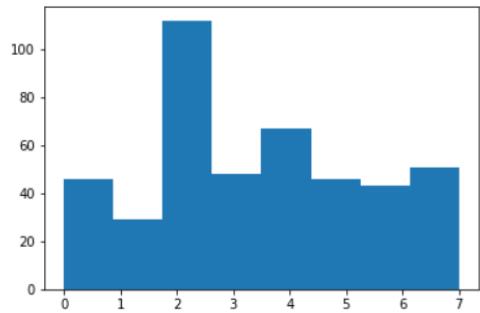


FIGURE 4.24: Histogram for Spectral Clustering (after PCA, without outliers)

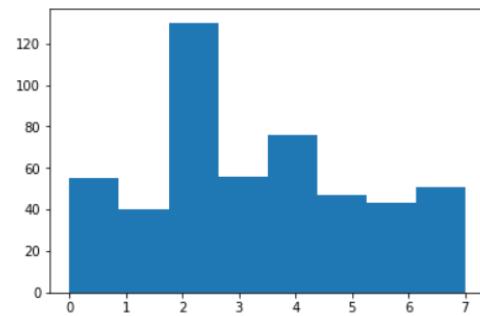


FIGURE 4.25: Histogram for Spectral Clustering (Placing back outliers)

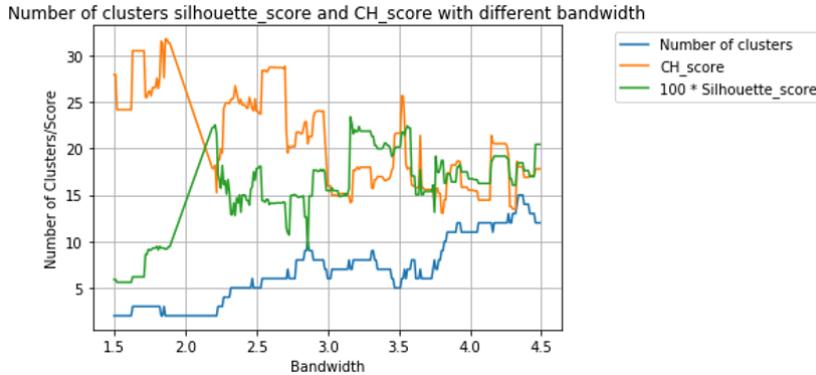


FIGURE 4.26: Mean Shift Algorithm scores(before PCA, with outliers)

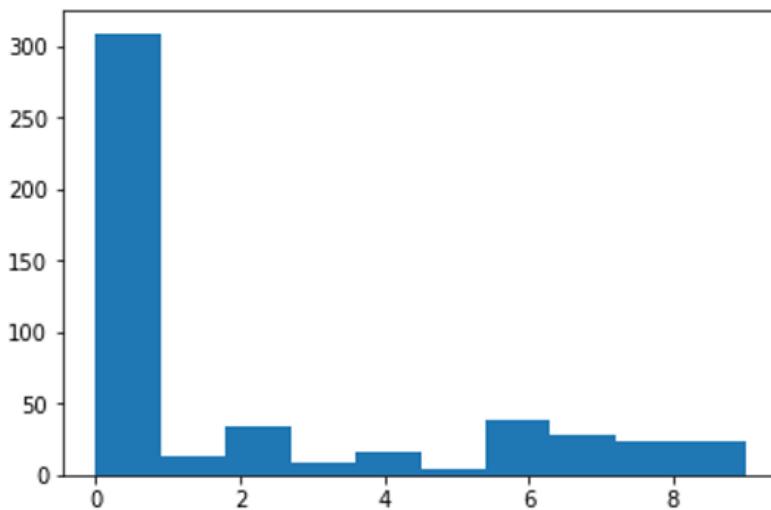


FIGURE 4.27: Histogram for Mean Shift Clustering(before PCA, with outliers)

4.6 Mean Shift

As mentioned in the background part, Mean Shift Algorithm has so many merits that make it a robust model. We expect the results to be good and let's try our four datasets. What's more, there is only one hyper parameter needed to be tuned (bandwidth).

4.6.1 Dataset before PCA with outliers

As presented in the figure 4.26, the results of Mean Shift model look fine when the bandwidth goes above 2.5: The number of clusters are around 10 and the CH scores are acceptable as well.

However, when we focus on the number of the samples in different clusters, it is imbalanced as shown in Figure 4.27.

This result can serve as benchmark and we can compare the result of this dataset with those of the other three to test whether PCA and Outliers detection enhance model performance.

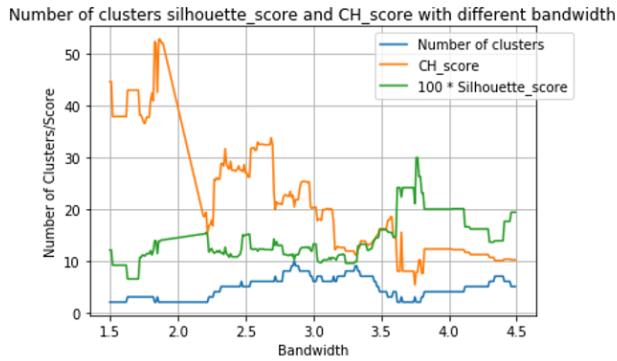


FIGURE 4.28: Mean Shift Algorithm scores(before PCA, without outliers)

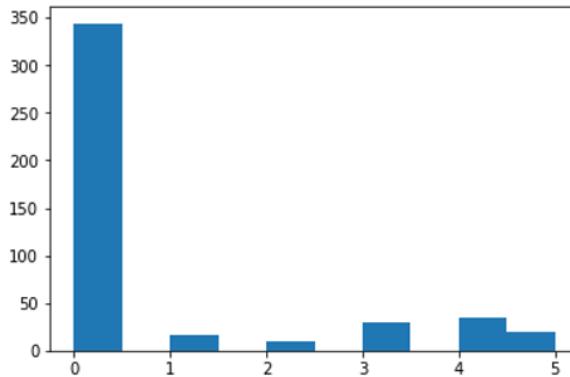


FIGURE 4.29: Histogram for Mean Shift Clustering Result(before PCA, without outliers)

4.6.2 Dataset before PCA without outliers

The result of this dataset is similar to last dataset. Dealing with outliers doesn't make the model performance better.(Figure 4.28)

No matter what value we pick for bandwidth, we always end up with imbalanced clusters.(Figure 4.29)

We can come to conclusion that removing outliers doesn't help us improve Mean Shift model with features before PCA.

4.6.3 Dataset after PCA with outliers and without outliers

Similar situation happens in these two datasets. Although CH score rises compared to model without PCA. The results are still imbalanced.(Figure 4.30, 4.31)

4.6.4 Model Conclusion

In all, Mean Shift model doesn't work well on the features no matter if we use PCA or Outliers Removal method.

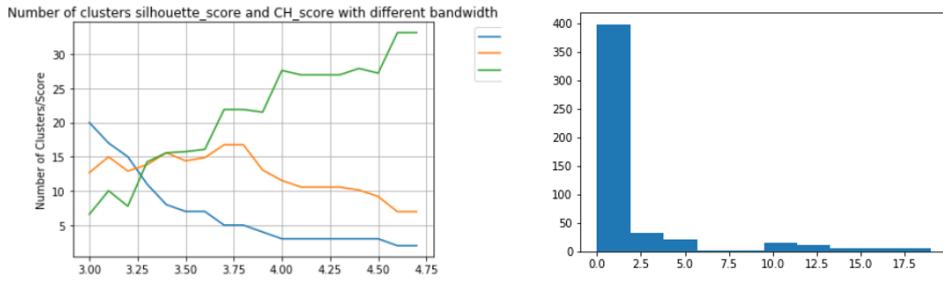


FIGURE 4.30: Mean Shift Clustering Result (After PCA, with outliers)

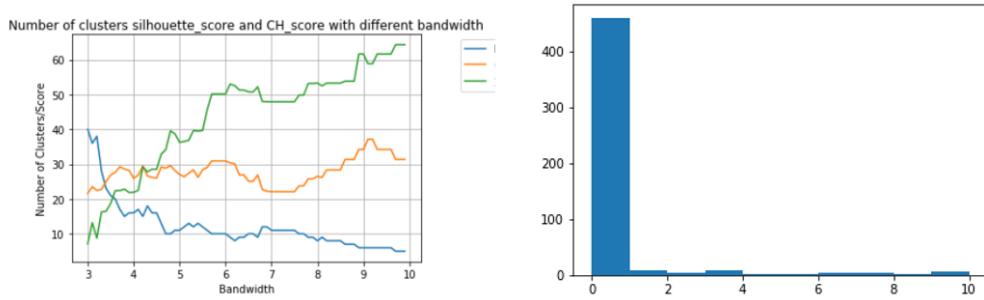


FIGURE 4.31: Mean Shift Clustering Result (After PCA, without outliers)

4.7 DBSCAN

DBSCAN is derived from spatial clustering and it is density-based and able to deal with noise. However, the results are similar to spatial clustering. DBSCAN performs badly for all distance matrixes. The major problem is imbalanced clusters as well.

For distance matrix, Cityblock, Cosine, Euclidean, L1, L2 and Manhattan are tested. And there are two hyper parameters need to be tuned. Grid Search is used and we filter those results have around 10 clusters and relatively high CH score.

However, the distributions of the number of the samples in clusters are imbalanced shown in the figure 4.32.

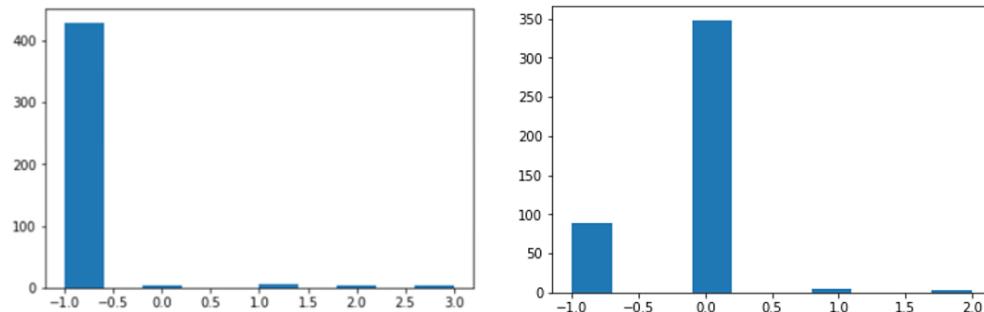


FIGURE 4.32: DBSCAN Clustering Result

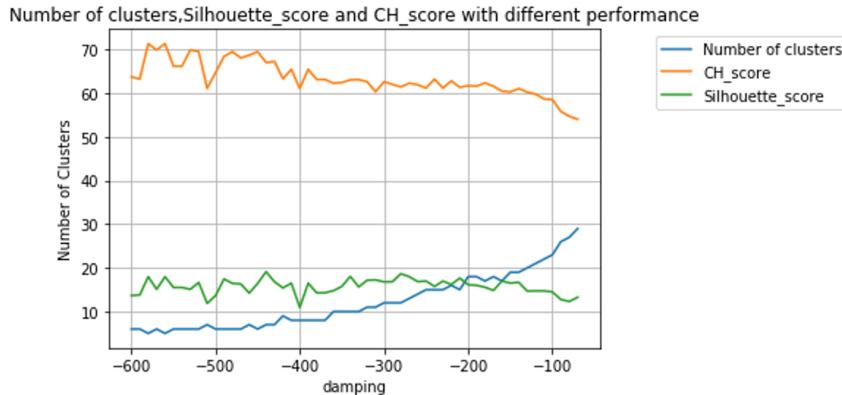


FIGURE 4.33: Affinity Propagation Clustering Scores (Before PCA, with outliers)

DBSCAN requires the density of the samples invariant within each cluster. This assumption might not hold for either of our four datasets. That could be reason DBSCAN does not work.

4.8 Affinity Propagation

Preference is a crucial hyper parameter in Affinity Propagation. The number of clusters varies as preference changes.

4.8.1 Dataset before PCA with outliers

We still start with dataset before PCA with outliers. The results are much better than those of DBSCAN and Mean Shift. CH score is around 60 and the outcomes are balanced. (Figure 4.33)

However, the only drawback is that some of clusters only have one or two samples. The outliers have negative effects on the Affinity Propagation algorithm.(Figure 4.34)

This dataset without PCA and outliers can serve as benchmark to the other three datasets. We can see how PCA and Outliers Detection effect the results.

4.8.2 Dataset before PCA without outliers

Classify samples without outliers

We remove outliers according to 5-sigma rule and train the model. Similar to the benchmark, we also get balanced results with high CH score. (Figure 4.35)

One improvement is that there is not any cluster only with one or two samples. All clusters contain at least 10 samples, which means our model is not creating clusters just for the outliers.(Figure 4.36)

We picked -150 as the final value for hyper parameter, preference, considering the number of clusters and CH score. The CH score is 43 and the number of clusters is 9 for this case.

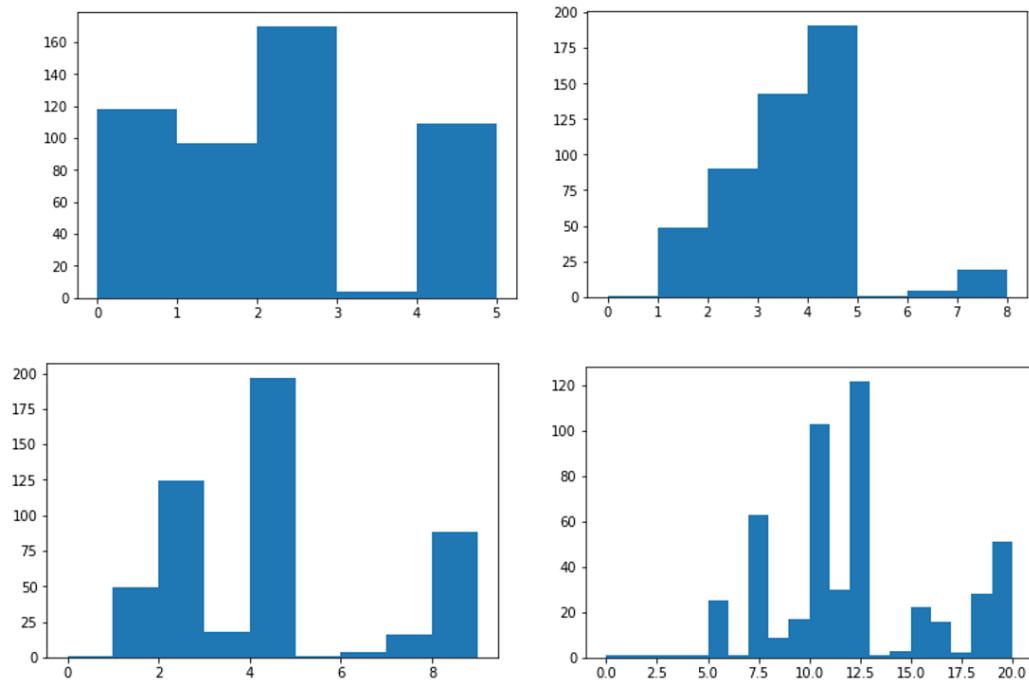


FIGURE 4.34: Histogram for Affinity Propagation Clustering Result
(Before PCA with outliers)

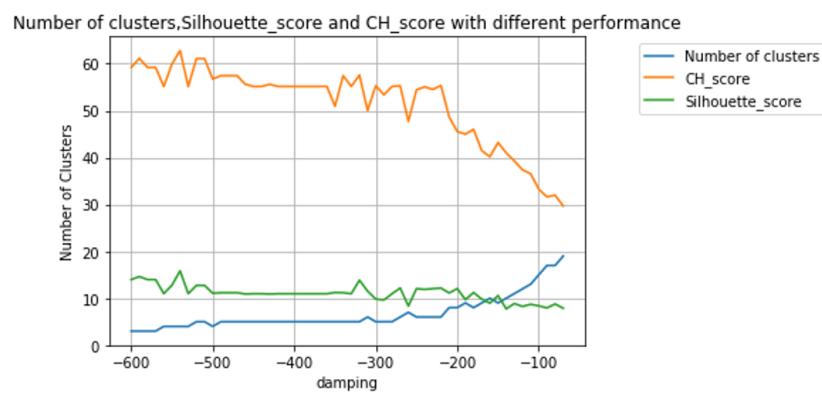


FIGURE 4.35: Affinity Propagation Clustering Scores (Before PCA,
without outliers)

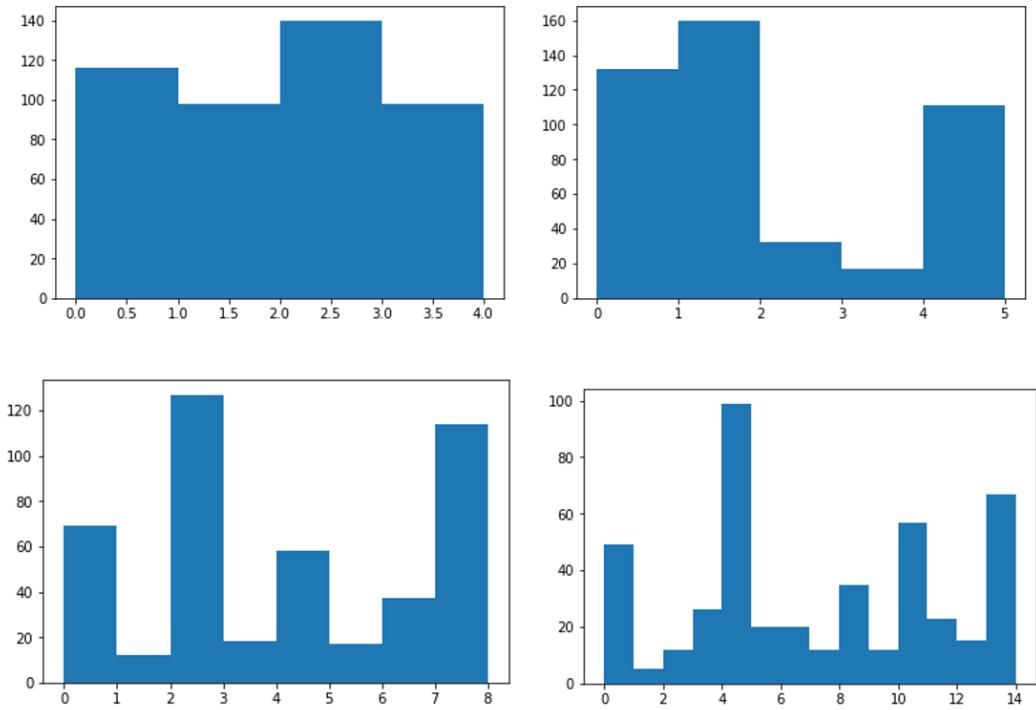


FIGURE 4.36: Histogram for Affinity Propagation Clustering Result
(Before PCA without outliers)

Assign outliers back to clusters

Then we assign the outliers back to clusters according to the distance between them and centers of the clusters. For each outlier, we calculate the distances between outlier and each clustering center. Then find the shortest distance and assign the outlier to that clustering center.

After assigning the outliers back to clusters. CH score decreases significantly, from 43 to 31, which indicates that the outliers do not fit in any clusters.

4.8.3 Dataset after PCA with outliers

After PCA transformation, CH score remain unchanged compared with benchmark.(Figure 4.37)

And for the distribution of the samples, there are 5 clusters that only contain less than 3 samples. These clusters are true clusters, just holders for outliers. (Figure 4.38)

4.8.4 Dataset after PCA with outliers

Classify samples without outliers

After PCA and Outliers removal, the CH score is much better than the benchmark. PCA and Outliers removal bring significant positive effects to the model.(Figure 4.39)

And as shown in the figure 4.40, the number of samples in each cluster is balanced and there are no fake cluster contain only one or two sample. These are quite satisfying result.

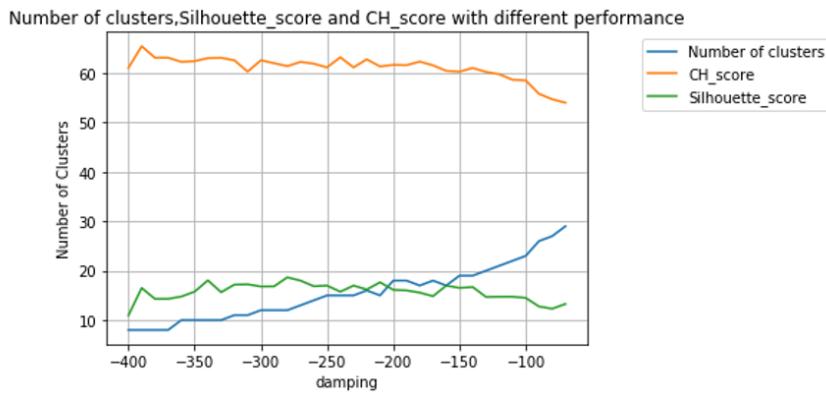


FIGURE 4.37: Affinity Propagation Clustering Scores (After PCA, with outliers)

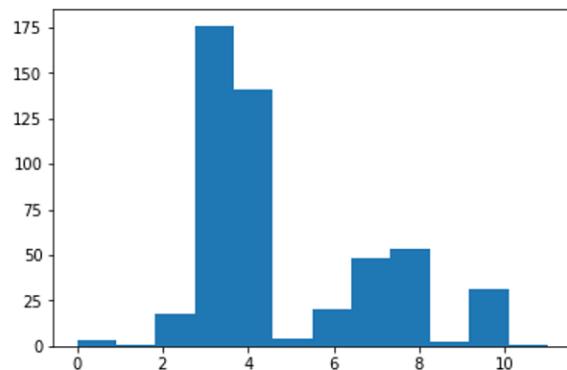


FIGURE 4.38: Histogram for Affinity Propagation Clustering Result (After PCA without outliers)

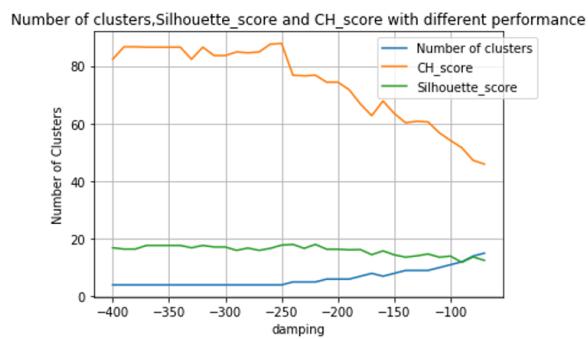


FIGURE 4.39: Affinity Propagation Clustering Scores (After PCA, without outliers)

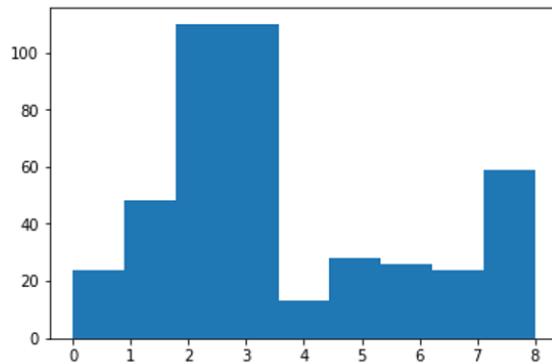


FIGURE 4.40: Histogram for Affinity Propagation Clustering Result
(After PCA without outliers)

Assign outliers back to clusters

We assign outliers back to clusters, but CH score decrease by 20, from 80 to 60.

4.8.5 Model Conclusion

To sum up, Affinity Propagation work well in these datasets and beats DBSCAN and Mean Shift. Outliers removal can improve the model performance and can eliminate fake clusters that only have one or two samples. Among the four datasets, the one with both PCA and outliers removal work the best: balanced result with 80 CH score.

Chapter 5

Model Comparison

5.1 The distribution of the results

According to the evenness of the results of the models, the clustering methods can be classified into two groups:

- 1) The models with balanced results: Affinity Propagation, K-Mean, MiniBatch K-Mean, Birch, Spectral and Agglomerative.
- 2) The models with imbalanced results: Mean Shift, DBSCAN.

For those models that have reasonable results, we try to rank them by using Silhouette Coefficient and CH score separately.

5.2 Silhouette Coefficient and CH score Comparisons

In order to have a better idea of the results of the models, we rank the best model of each clustering method according to silhouette scores. Best models should be models that have balanced results and around ten clusters. We show the ranking in the Table 5.1.

The differences among models Silhouette Scores are not that significant. However, we can't jump into conclusion that Agglomerative is the best model and Birch is the worst among these models.

Then we turn to see CH scores of different models.(Table 5.2)

Combining CH score and Silhouette score, we can say that Birch works the worst in these models. Agglomerative clustering and Spectral clustering work better than others.

TABLE 5.1: Model Silhouette Score Comparison

Clustering Methods	Silhouette Score
Agglomerative	0.21
Kmean	0.18
Spectral	0.16
MiniBatch-Kmean	0.15
Affinity Propagation	0.15
Birch	0.13

TABLE 5.2: Model CH Score Comparison

Clustering Methods	CH Score
Spectral	69
Affinity Propagation	61
Agglomerative	60
Kmean	58
MiniBatch-Kmean	53
Birch	51

5.3 PCA transformation effects and outliers removal effects

For features after PCA, each model gets higher CH score and Silhouette score, usually 10% to 30% higher. This convinces us that PCA transformation can help clustering models to correctly classify the samples.

For outliers removal, it is also proven a necessary process in data processing. If dataset with outliers is used, some clusters with only a few samples (under 5 samples) will be created for outliers. These clusters only contain some samples that are so special that can't be assigned to other clusters. However, after outliers detection and removal, there are no fake clusters created anymore. Every cluster has over 20 samples.

In all, PCA transformation and outliers removal are beneficial methods that can improve models performance and clustering results.

5.4 Outliers Reassignment

After we finish clustering the dataset without outliers, we assign the outliers we removed back to the clusters and we have a look at CH score and silhouette score. As expected, they drop accordingly. However, the drop is acceptable. Outliers certainly will have negative effects on the scores because they do not belong to any of the existing clusters. The reason we still want to assign them back is that If we have to give a label to these abnormal samples, we still can give them labels and the process of giving them labels will not effect we classify those normal samples.

Chapter 6

Conclusion and Future Research

In our capstone project, data is retrieved from Yahoo Finance. We clean the data, deal with outliers and do the feature engineering. After that, PCA analysis and transformation are used to perform dimensions reduction. Then, data is fed into different clustering models. After we tune the hyper parameters in the models, we find out which models can classify the samples correctly without imbalanced issue. Also, Evaluation with CH score and silhouette score are also hired to rank the models. There are still works left needed to be done. What's more, we manage to figure out how PCA and outliers removal improve our models' performance.

However, there are still works left to be done. There are still room for us to improve the features. More work about the features engineering can be done. Since we only have 6 features for technical features, more technical features can be tested and added to framework. Besides, the number of fundamental features is much greater than that of technical features. This makes fundamental features more important than technical features when we do the training. Weights can be assigned to make them equally important.

What's more, Equivalency test and Distance test can be performed in future research. Equivalency test can assist to find the most similar products within a universe. Distance analysis can be utilized to determine the effectiveness of the equivalency test. In addition, we can combine the results of previous clustering methods and try to test the consistency of the results. Question like are there some pairs of stocks always been clustered into one group can be solved.

Appendix A

Code: Data Preprocessing

A.1 Code: Get Data

```

1 # regular imports
2 import pandas as pd
3 import numpy as np
4
5 # for stock features import
6 import yahoo_historical
7 import yahoofinancials
8 from yahoofinancials import YahooFinancials
9
10 # get ticker table
11 ticker_url = "https://en.wikipedia.org/wiki/List_of_S%26P_500_companies"
12
13 ticker_table = pd.read_html(ticker_url)
14
15 # generate tickers list
16 tickers = ticker_table[0][0][1: ].tolist()
17
18 # get technical data
19 # get the historical price data
20 sp_data = YahooFinancials(tickers)
21 sp_prices = sp_data.get_historical_price_data('2018-06-01', '2019-06-01'
22     , 'daily')
23
24 result = []
25 for stock in tickers:
26     # check whether we got the prices data for a specific stock
27     # if not, pass to the next
28     if "prices" not in sp_prices[stock]:
29         continue
30
31     temp_list = sp_prices[stock]["prices"]
32
33     # reformatting the list
34     for daily_data in temp_list:
35         result_temp = []
36         result_temp.append(stock)
37         result_temp.append(daily_data["formatted_date"])
38         result_temp.append(daily_data["high"])
39         result_temp.append(daily_data["low"])
40         result_temp.append(daily_data["adjclose"])
41         result_temp.append(daily_data["volume"])
42         result.append(result_temp)
43
44 # store the data into pandas dataframe
45 sp_df = pd.DataFrame(np.array(result))
46 sp_df.columns = ["ticker", "date", "high", "low", "adjclose", "volume"]

```

```
45 # output to csv file
46 sp_df.to_csv(r'sp500.csv')
47
48
49 # get fundamental data
50 # get the key statistics of the stocks
51 sp_fundamental = sp_data.get_key_statistics_data
52
53 df_fundamental = pd.DataFrame.from_dict(sp_fundamental, orient = 'index'
54     )
55
56 sp_eps = sp_data.get_earnings_per_share()
57 sp_mc = sp_data.get_market_cap()
58 sp_pe = sp_data.get_pe_ratio()
59 sp_ni = sp_data.get_net_income()
60 sp_ebit = sp_data.get_ebit()
61
62 df_eps = pd.DataFrame(sp_eps, index=['eps']).T
63 df_mc = pd.DataFrame(sp_mc, index=['Market Cap']).T
64 df_pe = pd.DataFrame(sp_pe, index=['PE']).T
65 df_ni = pd.DataFrame(sp_ni, index=['NI']).T
66 df_ebit = pd.DataFrame(sp_ebit, index=['EBIT']).T
67
68 df = df_fundamental.join(df_eps, how = 'outer')
69 df = df.join(df_mc, how = 'outer')
70 df = df.join(df_pe, how = 'outer')
71 df = df.join(df_ni, how = 'outer')
72 df = df.join(df_ebit, how = 'outer')
73
74 # output to csv file
75 df.to_csv(r'sp500_fundamental_eps_mc.csv')
```

A.2 Code: Feature Engineering

```

1 import pandas as pd
2 import numpy as np
3
4 from sklearn.pipeline import Pipeline
5 try:
6     from sklearn.impute import SimpleImputer # Scikit-Learn 0.20+
7 except ImportError:
8     from sklearn.preprocessing import Imputer as SimpleImputer
9
10 from sklearn.preprocessing import StandardScaler
11 from sklearn.base import BaseEstimator, TransformerMixin
12 from scipy.stats import skew, mstats
13 from scipy.stats import kurtosis
14
15 import seaborn as sns
16 import matplotlib.pyplot as plt
17
18 from sklearn.decomposition import PCA
19
20 data_f = pd.read_csv(r"Fundamental_features.csv", index_col=0)
21 # f for Fundamental
22 data_t = pd.read_csv(r"Tech_Features.csv", index_col = 0)
23 # t for Tech
24
25 # Get SP500 index from yahoo finance for feature engineering purpose.
26
27 from yahoofinancials import YahooFinancials
28 yahoo_financials = YahooFinancials("^GSPC")
29 historical_stock_prices = yahoo_financials.get_historical_price_data(
    '2018-06-01', '2019-06-01', 'daily')
30 sp500 = pd.DataFrame(historical_stock_prices["^GSPC"]["prices"])
31
32
33 # Impute Missing Values and Feature Engineering
34
35 # technical features
36 # find ticker short of dates
37 ticker_delete = []
38 tickers = data_t["ticker"].unique()
39 rst = data_t.groupby('ticker')['date'].count()
40 for ticker in tickers:
41     if rst[ticker] != int(rst.median()):
42         ticker_delete.append(ticker)
43
44 print("Stock with unmatching data:")
45 print(ticker_delete)
46 for ticker in ticker_delete:
47     data_t = data_t.drop(data_t[data_t["ticker"]== ticker ].index, axis
        = 0)
48
49 # get close price, volume and return dataframes
50 close_df = data_t.pivot(index='ticker', columns='date', values='
    close_price')
51 volume_df = data_t.pivot(index='ticker', columns='date', values='volume'
    )
52 return_df = (close_df.pct_change(axis='columns')).drop(close_df.columns
    [0], axis=1)
53
54 # feature engineering for technical features
55 return_sp = ((sp500.adjclose.diff()/sp500.adjclose)[1:]).mean()

```

```

56 return_std = return_df.std(axis=1)
57 return_skew = return_df.skew(axis=1)
58 return_kurtosis = return_df.kurtosis(axis=1)
59 Max_Drawdown = (1 - close_df/close_df.cummax(axis = 1)).max(axis=1)
60 sharpe_ratio = (return_df.mean(axis=1)-return_sp)/return_std
61 volume_std = volume_df.std(axis=1)
62 t_feature = pd.DataFrame()
63 t_feature = t_feature.assign(Return_std = return_std, Return_skew =
64     return_skew, Return_kurtosis = return_kurtosis, Sharpe_Ratio =
65     sharpe_ratio, Volume_std = volume_std, Maxdrawdown = Max_Drawdown)
66
67 # fundamental features
68 # drop empty columns
69 data_f = data_f.dropna(axis='columns', how='all')
70 data_f = data_f[:-2]
71
72 # drop date columns and useless columns
73 data_f = data_f.drop(['fundInceptionDate', 'sharesShortPreviousMonthDate',
74     'lastSplitDate', 'lastSplitFactor', 'priceHint', 'dateShortInterest',
75     'maxAge', 'mostRecentQuarter', 'lastFiscalYearEnd', 'nextFiscalYearEnd'],
76     axis=1)
77
78 # copy to data_f
79 f_feature = data_f.copy()
80
81 # reset the index
82 f_feature.index = data_f["Ticker"]
83 f_feature = f_feature.drop(["Ticker"],axis = 1)
84
85 # in case there are infinity entries
86 f_feature = f_feature.replace([np.inf, -np.inf,"infinity","Infinity"],
87     np.nan)
88
89 # Feature Combination
90 # show the difference between two dataframe
91 t_inx = t_feature.index
92 f_inx = f_feature.index
93 display(t_inx.difference(f_inx))
94 display(f_inx.difference(t_inx))
95
96 feature_df = t_feature.join(f_feature, how='inner')
97
98 # Impute Missing Value and Standardization using Pipeline
99
100 class DataFrameSelector(BaseEstimator, TransformerMixin):
101     def __init__(self, attribute_names):
102         self.attribute_names = attribute_names
103     def fit(self, X, y=None):
104         return self
105     def transform(self, X):
106         return X[self.attribute_names].values
107
108 features = feature_df.columns
109
110 Preprocess_pipeline = Pipeline([
111     ("select_numeric", DataFrameSelector( features )), # set the
112     parameter for DataDrameSelector
113     ("imputer", SimpleImputer(strategy="mean")),
114     ("StandardScaler", StandardScaler(copy=True, with_mean=True,
115     with_std=True))
116 ])

```

```
111  
112 train_array = Preprocess_pipeline.fit_transform(feature_df)  
113  
114 train_df = pd.DataFrame(train_array)  
115 train_df.index = feature_df.index  
116 train_df.columns = feature_df.columns  
117 train_df.head()
```

A.3 Code: Principal Components Analysis and Outlier Processing

```

1 # continue with Code: Feature Engineering
2
3 # Principal Components Analysis
4
5 def heatMap(df, mirror):
6
7     # Create Correlation df
8     corr = df.corr()
9     # Plot figsize
10    fig, ax = plt.subplots(figsize=(20, 20))
11    # Generate Color Map
12    colormap = sns.diverging_palette(220, 10, as_cmap=True)
13
14    if mirror == True:
15        # Generate Heat Map, allow annotations and place floats in map
16        sns.heatmap(corr, cmap=colormap, annot=True, fmt=".2f")
17        # Apply xticks
18        plt.xticks(range(len(corr.columns)), corr.columns);
19        # Apply yticks
20        plt.yticks(range(len(corr.columns)), corr.columns)
21        # show plot
22    else:
23        # Drop self-correlations
24        dropSelf = np.zeros_like(corr)
25        dropSelf[np.triu_indices_from(dropSelf)] = True
26        # Generate Color Map
27        colormap = sns.diverging_palette(220, 10, as_cmap=True)
28        # Generate Heat Map, allow annotations and place floats in map
29        sns.heatmap(corr, cmap=colormap, annot=True, fmt=".1f", mask=
30 dropSelf)
31        # Apply xticks
32        plt.xticks(range(len(corr.columns)), corr.columns);
33        # Apply yticks
34        plt.yticks(range(len(corr.columns)), corr.columns)
35    # show plot
36    plt.show()
37
38 # draw heat map
39 heatMap(train_df, False)
40
41 pca_feature = PCA(random_state=100)
42 _ = pca_feature.fit(train_df)
43
44 acc_va_ratio = np.cumsum(pca_feature.explained_variance_ratio_)
45 # output variance ratio
46 acc_va_ratio
47
48 # draw variance ratio plot
49 PCA_var_ratio=plt.figure()
50 plt.plot(acc_va_ratio)
51 plt.legend(["Accumulated Variance Ratio"], bbox_to_anchor=(1.1,1))
52 plt.grid()
53 plt.title("Accumulated Variance Ratio")
54 plt.xlabel('Components')
55 plt.ylabel('Percentage')
56 plt.axhline(y=0.8)
57 _ = plt.show()

```

```
58 # set threshold to 0.8 and keep first 11 components
59 m_pca = PCA(n_components=11).fit_transform(train_df)
60
61
62 # Dealing with Outliers
63 # before PCA
64 # Without pca, samples with outliers:
65 outliers_tickers = np.array(train_df[(train_df>5)|(train_df<-5)].dropna(
66     how = 'all',axis = 0).index)
67 data_o = train_df.loc[outliers_tickers,:]
68 # Without pca, samples without outliers
69 data_wo_o = train_df[(train_df<=5)&(train_df>=-5)].dropna(axis = 0)# 4
70     clustering method
71
72 # after PCA
73 # create dataframe for features after pca
74 pca_df = pd.DataFrame(m_pca)
75 pca_df.columns = [ "{}th component".format(i+1) for i in range(11)]
76 pca_df.index = train_df.index
77
78 # outliers out of range of(-5,5)
79 # after pca ,samples outliers:
80 outliers_tickers = np.array(pca_df[(pca_df>5)|(pca_df<-5)].dropna(how =
81     'all',axis = 0).index)
82 pca_o = pca_df.loc[outliers_tickers,:]
83 # after pca,samples without outliers
84 pca_wo_o = pca_df[(pca_df<=5)&(pca_df>=-5)].dropna(axis = 0)
```


Appendix B

Code: Model Implementation

B.1 Model: KMeans

```

1  from sklearn.cluster import KMeans
2  from sklearn.cluster import MiniBatchKMeans
3  from sklearn.cluster import Birch
4  from sklearn.metrics import silhouette_samples, silhouette_score
5  from sklearn.metrics import calinski_harabaz_score
6
7
8 # Evaluation method with sh score
9
10 def find_optimal_clusters1(n,data):
11     cluster_num = np.arange(2,n)
12     scores = []
13     for k in cluster_num:
14         kmean = KMeans(n_clusters = k,random_state = 42,algorithm='full')
15         kmean.fit(data)
16         labels = kmean.predict(data)
17         score = silhouette_score(data,labels)
18         scores.append(score)
19
20     print("Optimal number for clusters:{}".format(cluster_num[np.argmax(scores)]))
21     print("Optimal SH score:{}".format(max(scores)))
22     MS_cluster=plt.figure()
23     plt.plot(cluster_num,scores)
24     plt.legend(["Number of clusters","SH_score"],bbox_to_anchor=(1.1,1))
25     plt.grid()
26     plt.title("Number of clusters and SH_score of Kmean")
27     plt.xlabel('Number of Clusters')
28     plt.ylabel('Score')
29     _ = plt.show()
30
31 # Modeling with After PCA Dataframe
32 find_optimal_clusters1(15,pca_wo_o) # max clustring range as 15
33
34 find_optimal_clusters1(15,data_wo_o)
35
36 # Elbow Method
37 def find_optimal_clusters2(n,data):
38     cluster_num = np.arange(2,n)
39     scores = []
40     for k in cluster_num:
41         kmean = KMeans(n_clusters = k, random_state = 42,algorithm='full').fit(data)
42         labels = kmean.predict(data)

```

```

43     score = kmean.inertia_
44     scores.append(score)
45
46
47 MS_cluster=plt.figure()
48 plt.plot(cluster_num,scores,'o-')
49 plt.legend(["Number of clusters"],bbox_to_anchor=(1.1,1))
50 plt.grid()
51 plt.title("Number of clusters and Sum of square of Kmean")
52 plt.xlabel('Number of Clusters')
53 plt.ylabel('Score')
54 _ = plt.show()
55
56 # With optimal parameters and after PCA dataframe.
57 find_optimal_clusters2(15,pca_wo_o)
58
59
60 # Gap statistic method
61 def find_optimal_clusters3(data, nrefs=3, maxClusters=15):
62     """
63         Calculates KMeans optimal K using Gap Statistic from Tibshirani,
64         Walther, Hastie
65         Params:
66             data: ndarray of shape (n_samples, n_features)
67             nrefs: number of sample reference datasets to create
68             maxClusters: Maximum number of clusters to test for
69         Returns: (gaps, optimalK)
70     """
71     gaps = np.zeros((len(range(1, maxClusters)),))
72     resultsdf = pd.DataFrame({'clusterCount':[], 'gap':[]})
73     for gap_index, k in enumerate(range(1, maxClusters)):
74
75         # Holder for reference dispersion results
76         refDisps = np.zeros(nrefs)
77
78         # For n references, generate random sample and perform kmeans
79         # getting resulting dispersion of each loop
80         for i in range(nrefs):
81
82             # Create new random reference set
83             randomReference = np.random.random_sample(size=data.shape)
84
85             # Fit to it
86             km = KMeans(k)
87             km.fit(randomReference)
88
89             refDisp = km.inertia_
90             refDisps[i] = refDisp
91
92             # Fit cluster to original data and create dispersion
93             km = KMeans(k)
94             km.fit(data)
95
96             origDisp = km.inertia_
97             origDisp = km.inertia_
98
99             # Calculate gap statistic
100            gap = np.log(np.mean(refDisps)) - np.log(origDisp)
101
102            # Assign this loop's gap statistic to gaps
103            gaps[gap_index] = gap

```

```

103     resultsdf = resultsdf.append({'clusterCount':k, 'gap':gap},
104     ignore_index=True)
105
106     return (gaps.argmax() + 1, resultsdf) # Plus 1 because index of 0
107     means 1 cluster is optimal, index 2 = 3 clusters are optimal
108
109 # Use gap statistic to evaluate kmeans model
110 k, gapdf = find_optimal_clusters3(pca_wo_o, nrefs=5, maxClusters=15)
111 print ('Optimal k is: ', k)
112 plt.plot(gapdf.clusterCount, gapdf.gap, linewidth=3)
113 plt.scatter(gapdf[gapdf.clusterCount == k].clusterCount, gapdf[gapdf.
114     clusterCount == k].gap, s=250, c='r')
115 plt.grid(True)
116 plt.xlabel('Cluster Count')
117 plt.ylabel('Gap Value')
118 plt.title('Gap Values by Cluster Count')
119 plt.show()
120
121 # Prediction distribution
122 plt.figure(figsize=(14, 6))
123 plt.subplot(121)
124 plt.hist(KMeans(8,random_state = 42,algorithm='full').fit_predict(
125     pca_wo_o))
126 plt.title("Distribution with 8 Clustering group")
127 plt.subplot(122)
128 plt.hist(KMeans(11,random_state = 42,algorithm='full').fit_predict(
129     pca_wo_o))
130 plt.title("Distribution with 11 Clustering group")
131
132 # outlier assignment
133 from scipy.spatial import distance
134 # center can be list or array
135 # assume each center correponds to label of [0,1,2...]
136 # return a list of label
137 def classify_outliers(centers,outliers_df):
138     labels = []
139     for sample in np.array(outliers_df):
140         min_center = float('inf')
141         min_center_th = -10086
142         for i,center in enumerate(centers):
143             dst = distance.euclidean(center,sample)
144             if dst<min_center:
145                 min_center = dst
146                 min_center_th = i
147         labels.append(min_center_th)
148     return labels
149
150 # example of using this fuction
151 # centers are calculated from the models, lets assume two random
152 # centers here
153 kmeans_11 = KMeans(11,random_state = 42,algorithm='full').fit(pca_wo_o)
154 centers = kmeans_11.cluster_centers_
155 outlier_labels = classify_outliers(centers,pca_o)
156 pca_wo_o_labels = kmeans_11.labels_
157 labels = np.append(pca_wo_o_labels,outlier_labels)
158
159 score = silhouette_score(pca_df,labels)
160 score2 = silhouette_score(pca_wo_o,pca_wo_o_labels)
161 print(score)
162 print(score2)

```

B.2 Model: MiniBatchKMeans

```

1 def find_optimal_clusters_mini(n,data,batch_size):
2     cluster_num = np.arange(2,n)
3     scores = []
4     for k in cluster_num:
5         minikmean = MiniBatchKMeans(n_clusters = k,random_state = 42,
6 batch_size = batch_size).fit(data)
7         labels = minikmean.predict(data)
8         score = silhouette_score(data,labels)
9         scores.append(score)
10
11     print("Optimal number for clusters:{}".format(cluster_num[np.argmax
12 (scores)]))
13     print("Optimal SH score:{}".format(max(scores)))
14     MS_cluster=plt.figure()
15     plt.plot(cluster_num,scores)
16     plt.legend(["Number of clusters","SH_score"],bbox_to_anchor=(1.1,1)
17 )
18     plt.grid()
19     plt.title("Number of clusters and SH_score of MiniBatchKmean")
20     plt.xlabel('Number of Clusters')
21     plt.ylabel('CH Score')
22     _ = plt.show()
23
24
25 # Elbow Method
26 def find_optimal_clusters_mini2(n,data,batch_size):
27     cluster_num = np.arange(2,n)
28     scores = []
29     for k in cluster_num:
30         minikmean = MiniBatchKMeans(n_clusters = k,random_state = 42,
31 batch_size = batch_size).fit(data)
32         labels = minikmean.predict(data)
33         score = minikmean.inertia_
34         scores.append(score)
35
36     MS_cluster=plt.figure()
37     plt.plot(cluster_num,scores,'o-')
38     plt.legend(["Number of clusters"],bbox_to_anchor=(1.1,1))
39     plt.grid()
40     plt.title("Number of clusters and Sum of square of MiniBatchKmean")
41     plt.xlabel('Number of Clusters')
42     plt.ylabel('Score')
43     _ = plt.show()
44
45 find_optimal_clusters_mini2(15,pca_wo_o,10)
46
47
48 # Gap statistic method
49 def find_optimal_clusters_mini3(data, nrefs=3, maxClusters=15):
50     """
51     Calculates KMeans optimal K using Gap Statistic from Tibshirani,
52     Walther, Hastie
53     Params:
54         data: ndarray of shape (n_samples, n_features)
55         nrefs: number of sample reference datasets to create
56         maxClusters: Maximum number of clusters to test for

```

```

56     Returns: (gaps, optimalK)
57     """
58     gaps = np.zeros((len(range(1, maxClusters)),))
59     resultsdf = pd.DataFrame({'clusterCount':[], 'gap':[]})
60     for gap_index, k in enumerate(range(1, maxClusters)):
61
62         # Holder for reference dispersion results
63         refDisps = np.zeros(nrefs)
64
65         # For n references, generate random sample and perform kmeans
66         # getting resulting dispersion of each loop
67         for i in range(nrefs):
68
69             # Create new random reference set
70             randomReference = np.random.random_sample(size=data.shape)
71
72             # Fit to it
73             km = MiniBatchKMeans(k)
74             km.fit(randomReference)
75
76             refDisp = km.inertia_
77             refDisps[i] = refDisp
78
79             # Fit cluster to original data and create dispersion
80             km = MiniBatchKMeans(k)
81             km.fit(data)
82
83             origDisp = km.inertia_
84             origDisp = km.inertia_
85
86             # Calculate gap statistic
87             gap = np.log(np.mean(refDisps)) - np.log(origDisp)
88
89             # Assign this loop's gap statistic to gaps
90             gaps[gap_index] = gap
91
92             resultsdf = resultsdf.append({'clusterCount':k, 'gap':gap},
93                                         ignore_index=True)
94
95
96 k, gapdf = find_optimal_clusters_mini3(pca_wo_o, nrefs=5, maxClusters
97 =15)
98 print ('Optimal k is: ', k)
99 plt.plot(gapdf.clusterCount, gapdf.gap, linewidth=3)
100 plt.scatter(gapdf[gapdf.clusterCount == k].clusterCount, gapdf[gapdf.
101   clusterCount == k].gap, s=250, c='r')
102 plt.grid(True)
103 plt.xlabel('Cluster Count')
104 plt.ylabel('Gap Value')
105 plt.title('Gap Values by Cluster Count')
106 plt.show()
107
108 # Prediction distribution
109
110 plt.figure(figsize=(14, 6))
111 plt.subplot(121)
112 plt.hist(MiniBatchKMeans(n_clusters = 6,random_state = 42, batch_size =
113   10).fit_predict(pca_wo_o))
114 plt.title("Distribution with 8 Clustering group")

```

```
114
115 plt.subplot(122)
116 plt.hist(MiniBatchKMeans(n_clusters = 11,random_state = 42, batch_size
117 = 10).fit_predict(pca_wo_o))
118 plt.title("Distribution with 11 Clustering group")
119
120 # outlier assignment
121
122 #example of using this fuction
123 #centers are calculated from the models, lets assume two random centers
124 here
125 MiniBatchKMeans_6 = MiniBatchKMeans(n_clusters = 6,random_state = 42,
126 batch_size = 10).fit(pca_wo_o)
127 minicenters = MiniBatchKMeans_6.cluster_centers_
128 minioutlier_labels = classify_outliers(minicenters,pca_o)
129 minipca_wo_o_labels = MiniBatchKMeans_6.labels_
130 minilabels = np.append(minipca_wo_o_labels,minioutlier_labels)
131 miniscore = silhouette_score(pca_df,minilabels)
132 miniscore2 = silhouette_score(pca_wo_o,minipca_wo_o_labels)
133 print(miniscore)
134 print(miniscore2)
```

B.3 Model: Birch

```

1 # Grid search
2 def find_optimal_clusters_birch(data, bf, thd):
3     cluster_num = [2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14]
4     scores = []
5     for n in cluster_num:
6         brc = Birch(branching_factor=bf, n_clusters=n, threshold=thd,
7         compute_labels=True).fit(data)
8         label = brc.predict(data)
9         score = calinski_harabaz_score(data, label)
10        scores.append(score)
11    print("Optimal number for clusters:{}".format(cluster_num[np.argmax(
12        scores)]))
13    print("Optimal CH score:{}".format(max(scores)))
14
15    MS_cluster=plt.figure()
16    plt.plot(cluster_num,scores)
17    plt.legend(["Number of clusters","CH_score"],bbox_to_anchor=(1.1,1))
18    plt.grid()
19    plt.title("Number of clusters and CH_score of Birch")
20    plt.xlabel('Number of Clusters')
21    plt.ylabel('CH Score')
22    _ = plt.show()
23
24 def grid_search_birch(data):
25     threshold= [0.1,0.2,0.3,0.5,0.7,1]
26     branching_factor = [20, 30, 40, 50, 60, 80, 100]
27     for i in threshold:
28         for j in branching_factor:
29             try:
30                 find_optimal_clusters_birch(data,j,i)
31                 print(j,i)
32             except:
33                 print('not applicable')
34
35 grid_search_birch(pca_wo_o)
36
37 # find optimal cluster number
38 find_optimal_clusters_birch(pca_wo_o, 50 ,1)
39
40 # prediction distribution
41 bir = Birch(branching_factor=50, n_clusters=4, threshold=1,
42             compute_labels=True).fit(pca_wo_o)
43 bir2 = Birch(branching_factor=50, n_clusters=8, threshold=1,
44             compute_labels=True).fit(pca_wo_o)
45
46 plt.figure(figsize = (14,6))
47 plt.subplot(121)
48 plt.hist(bir.predict(pca_wo_o))
49 plt.title("Distribution with 4 Clustering group")
50
51 plt.subplot(122)
52 plt.hist(bir2.predict(pca_wo_o))
53 plt.title("Distribution with 8 Clustering group")
54
55 # outlier assignment
56 #example of using this fuction

```

```
55 #centers are calculated from the models, lets assume two random centers
56     here
57 bircenters = bir.subcluster_centers_
58 biroutlier_labels = classify_outliers(minicenters,pca_o)
59 birpca_wo_o_labels = bir.labels_
60 birlabels = np.append(birpca_wo_o_labels,biroutlier_labels)
61
62 birscore = silhouette_score(pca_df,birlabels)
63 birscore2 = silhouette_score(pca_wo_o,birpca_wo_o_labels)
64 print(birscore)
65 print(birscore2)
```

B.4 Model: Agglomerative

```

1  from scipy.cluster.hierarchy import dendrogram, linkage,
2      set_link_color_palette
3  from scipy.cluster.hierarchy import cophenet
4  from scipy.spatial.distance import pdist
5  import pylab
6
7  from sklearn.cluster import AgglomerativeClustering
8  from sklearn.metrics import calinski_harabaz_score, silhouette_score
9
10 # define grid search function for Agglomerative Clustering
11
12 def AC_grid_search(df,n,link_list,aff_list):
13     num_of_cluster = np.arange(2,n+1,1)
14     ch_score = []
15     sh_score = []
16     real_cluster = []
17     for k in num_of_cluster:
18         opt_ch = 0
19         opt_sh = 0
20         opt_link = 'null'
21         opt_aff = 'null'
22         n_cluster = 0
23         for i in link_list:
24             for j in aff_list:
25                 try:
26                     AC_model = AgglomerativeClustering(n_clusters=k,
27                     linkage = i, affinity=j)
28                     predicted = AC_model.fit_predict(df)
29                     temp_score = calinski_harabaz_score(X = df, labels =
30                     predicted )
31                     if temp_score > opt_ch:
32                         opt_ch = temp_score
33                         opt_sh = silhouette_score(X = df, labels =
34                     predicted)*100
35                     opt_link = i
36                     opt_aff = j
37                     n_cluster = len(set(predicted))
38                 except:
39                     pass
40                 print('For',k, 'clusters, Optimal combination is: linkage = ',
41                     opt_link, ',affinity = ', opt_aff)
42                 print('Highest CH_score is: ',opt_ch, 'Corresponding
43 SH_score is: ', opt_sh)
44                 ch_score.append(opt_ch)
45                 sh_score.append(opt_sh)
46                 real_cluster.append(n_cluster)
47
48                 AC_clusters=plt.figure()
49                 plt.plot(num_of_cluster,ch_score)
50                 plt.plot(num_of_cluster,sh_score)
51                 plt.plot(num_of_cluster, real_cluster)
52                 plt.grid()
53                 plt.legend(['ch_score', "sh_score * 100","real_cluster_num"],bbox_to_anchor=(1.1,1))
54                 plt.title("CH_score with different number of clusters")
55                 plt.xlabel('Number of Clusters')
56                 _ = plt.show()
57
58 # parameter list

```

```

54 affinity_list = ['euclidean','l1','l2','manhattan','cosine',''
55     'precomputed']
56 linkage_list = ['ward','complete', 'average', 'single']
57
58 # before PCA
59 AC_grid_search(train_df,20,linkage_list,affinity_list)
60
61 # Plot dendrogram
62 ward_nopca = linkage(train_df, 'ward')
63 plt.figure(figsize=(25, 10))
64 labelsize=20
65 ticksizes=15
66 set_link_color_palette(['firebrick','lightcoral','darkorange','gold',''
67     'yellowgreen','mediumaquamarine','lightseagreen','skyblue',''
68     'steelblue','slateblue','rebeccapurple','darkmagenta'])
69 plt.title('Agglomerative Clustering Dendrogram (ward, no PCA)','
70     fontsize=labelsizes)
71 plt.xlabel('stock', fontsize=labelsizes)
72 plt.ylabel('distance', fontsize=labelsizes)
73 dendrogram(
74     ward_nopca,
75     color_threshold = 28,
76     leaf_rotation=90., # rotates the x axis labels
77     leaf_font_size=8., # font size for the x axis labels
78 )
79 pylab.yticks(fontsize=ticksizes)
80 pylab.xticks(rotation=-90, fontsize=ticksizes)
81 plt.show()
82
83 # plot histogram for optimal result
84 ac_nopca2 = AgglomerativeClustering(n_clusters=12, linkage = 'ward')
85 ac_nopca_predicted2 = ac_nopca2.fit_predict(train_df)
86 print('ch_score = %.2f' % calinski_harabaz_score(X = train_df, labels =
87     ac_nopca_predicted2))
88 print('sh-score = %.2f' % silhouette_score(X = train_df, labels =
89     ac_nopca_predicted2))
90 plt.hist(ac_nopca_predicted2,bins=12)
91
92 # calculate std for clustering result
93 result = pd.DataFrame(ac_nopca_predicted2)
94 result[0].value_counts().std()
95
96
97 # after PCA, with outliers
98 AC_grid_search(m_pca,20,linkage_list,affinity_list)
99
100 # Dendrogram
101 ward_pca = linkage(m_pca, 'ward')
102 plt.figure(figsize=(25, 10))
103 labelsize=20
104 ticksizes=15
105 plt.title('Agglomerative Clustering Dendrogram (ward, PCA, with'
106     'outliers)', fontsize=labelsizes)
107 plt.xlabel('stock', fontsize=labelsizes)
108 plt.ylabel('distance', fontsize=labelsizes)
109 dendrogram(
110     ward_pca,
111     color_threshold = 28,
112     leaf_rotation=90., # rotates the x axis labels
113     leaf_font_size=8., # font size for the x axis labels
114 )
115 pylab.yticks(fontsize=ticksizes)
116 pylab.xticks(rotation=-90, fontsize=ticksizes)

```

```
110 plt.show()
111
112 # histogram
113 ac_pca1 = AgglomerativeClustering(n_clusters=9, linkage = 'ward')
114 ac_pca_predicted1 = ac_pca1.fit_predict(m_pca)
115 print('ch_score = %.2f' % calinski_harabaz_score(X = m_pca, labels =
116     ac_pca_predicted1))
117 print('sh-score = %.2f' % silhouette_score(X = m_pca, labels =
118     ac_pca_predicted1))
119 plt.hist(ac_pca_predicted1,bins=9)
120
121 # clustering result std
122 result = pd.DataFrame(ac_pca_predicted1)
123 result[0].value_counts().std()
124
125 # after PCA without outliers
126 AC_grid_search(pca_wo_o,20,linkage_list,affinity_list)
127
128 # Dendrogram
129 ward_pca = linkage(pca_wo_o, 'ward')
130 plt.figure(figsize=(25, 10))
131 labelsize=20
132 ticksize=15
133 plt.title('Agglomerative Clustering Dendrogram (ward, PCA, without
134     outliers)', fontsize=labelsize)
135 plt.xlabel('stock', fontsize=labelsize)
136 plt.ylabel('distance', fontsize=labelsize)
137 dendrogram(
138     ward_pca,
139     color_threshold = 21.5,
140     leaf_rotation=90., # rotates the x axis labels
141     leaf_font_size=8., # font size for the x axis labels
142 )
143 pylab.yticks(fontsize=ticksize)
144 pylab.xticks(rotation=-90, fontsize=ticksize)
145 plt.show()
146
147 # histogram
148 ac_wo_o1 = AgglomerativeClustering(n_clusters=6, linkage = 'ward')
149 ac_wo_o_predicted1 = ac_wo_o1.fit_predict(pca_wo_o)
150 print('ch_score = %.2f' % calinski_harabaz_score(X = pca_wo_o, labels =
151     ac_wo_o_predicted1))
152 print('sh-score = %.2f' % silhouette_score(X = pca_wo_o, labels =
153     ac_wo_o_predicted1))
154 plt.hist(ac_wo_o_predicted1,bins=6)
155
156 # result std
157 result = pd.DataFrame(ac_wo_o_predicted1)
158 result[0].value_counts().std()
159
160 # dendrogram for n = 13
161 ward_pca = linkage(pca_wo_o, 'ward')
162 plt.figure(figsize=(25, 10))
163 labelsize=20
164 ticksize=15
165 set_link_color_palette(['firebrick','lightcoral','darkorange','gold','
166     yellowgreen','limegreen','mediumaquamarine','lightseagreen','
167     skyblue','steelblue','slateblue','rebeccapurple','darkmagenta'])
168 plt.title('Agglomerative Clustering Dendrogram (ward, PCA, without
169     outliers)', fontsize=labelsize)
170 plt.xlabel('stock', fontsize=labelsize)
171 plt.ylabel('distance', fontsize=labelsize)
```

```

165 dendrogram(
166     ward_pca,
167     color_threshold = 12.4,
168     leaf_rotation=90., # rotates the x axis labels
169     leaf_font_size=8., # font size for the x axis labels
170 )
171 pylab.yticks(fontsize=ticksizes)
172 pylab.xticks(rotation=-90, fontsize=ticksizes)
173 plt.show()
174
175 # histogram for n = 13
176 ac_wo_o2 = AgglomerativeClustering(n_clusters=13, linkage = 'ward')
177 ac_wo_o_predicted2 = ac_wo_o2.fit_predict(pca_wo_o)
178 print('ch_score = %.2f' % calinski_harabaz_score(X = pca_wo_o, labels =
    ac_wo_o_predicted2))
179 print('sh-score = %.2f' % silhouette_score(X = pca_wo_o, labels =
    ac_wo_o_predicted2))
180 plt.hist(ac_wo_o_predicted2, bins=13)
181
182 # result std
183 result = pd.DataFrame(ac_wo_o_predicted2)
184 result[0].value_counts().std()
185
186 # labelling outliers
187 # define label function
188 from scipy.spatial import distance
189 def classify_outliers(center_df, outliers_df):
190     labels = []
191     for sample in np.array(outliers_df):
192         min_center = float('inf')
193         min_center_th = -10086
194         i=0
195         for center in np.array(center_df):
196             dst = distance.euclidean(center, sample)
197             if dst<min_center:
198                 min_center = dst
199                 min_center_th = i
200                 i=i+1
201         labels.append(min_center_th)
202     return labels
203
204 pca_wo_o1 = pca_wo_o.copy()
205 pca_wo_o1['label'] = ac_wo_o_predicted1
206
207 # calculate center for n = 6
208 ac_center1 = pd.DataFrame()
209 for i in range(6):
210     ac_center1[i] = np.mean(pca_wo_o1[pca_wo_o1['label']==i], axis=0)
211
212 ac_center1 = ac_center1.drop(['label']).transpose()
213
214 # plot histogram after labelling outliers
215 ac_o_labels1 = classify_outliers(ac_center1, pca_o)
216 print( calinski_harabaz_score(X =pd.concat([pca_wo_o, pca_o]) ,labels =
    list(ac_wo_o_predicted1)+ac_o_labels1 ))
217 plt.hist(list(ac_wo_o_predicted1)+ac_o_labels1, bins=6)
218
219 # result std
220 result = pd.DataFrame(ac_wo_o_predicted1)
221 result[0].value_counts().std()
222
223 # get clustering center for n =13
224 pca_wo_o2 = pca_wo_o.copy()

```

```
225 pca_wo_o2['label'] = ac_wo_o_predicted2
226
227 ac_center2 = pd.DataFrame()
228 for i in range(13):
229     ac_center2[i] = np.mean(pca_wo_o2[pca_wo_o2['label']==i],axis=0)
230
231 ac_center2 = ac_center2.drop(['label']).transpose()
232
233 # histogram
234 ac_o_labels2 = classify_outliers(ac_center2,pca_o)
235 print(calinski_harabaz_score(X =pd.concat([pca_wo_o,pca_o]) ,labels =
236     list(ac_wo_o_predicted2)+ac_o_labels2 ))
237 plt.hist(list(ac_wo_o_predicted2)+ac_o_labels2,bins=13)
238
239 # result std
240 result = pd.DataFrame(ac_wo_o_predicted2)
241 result[0].value_counts().std()
```

B.5 Model: Spectral

```

1 from sklearn.cluster import SpectralClustering
2
3 # define grid search for spectral clustering
4 def SC_grid_search(df,n,ninit,gamma,aff_list,neighbors):
5     num_of_cluster = np.arange(2,n+1,1)
6     ch_score = []
7     sh_score = []
8     real_cluster = []
9     for k in num_of_cluster:
10         opt_ch = 0
11         opt_sh = 0
12         opt_init = 0
13         opt_gamma = 0
14         opt_aff = 'null'
15         opt_nei = 0
16         n_cluster = 0
17         for i in ninit:
18             for j in gamma:
19                 for l in aff_list:
20                     for m in neighbors:
21                         try:
22                             SC_model = SpectralClustering(n_clusters=k,
23                             n_init= i, gamma= j, affinity=l, n_neighbors=m, random_state=42)
24                             predicted = SC_model.fit_predict(df)
25                             temp_score = calinski_harabaz_score(X = df,
26                             labels = predicted )
27                             if temp_score > opt_ch:
28                                 opt_ch = temp_score
29                                 opt_sh = silhouette_score(X = df,
30                                 labels = predicted)*100
31                                 opt_init = i
32                                 opt_gamma = j
33                                 opt_aff = l
34                                 opt_nei = m
35                                 n_cluster = len(set(predicted))
36                         except:
37                             pass
38                         print('For',k, 'clusters, Optimal combination is: n_init = ',
39                             opt_init, ', gamma = ', opt_gamma,
40                             ',affinity = ', opt_aff, ',number of neighbors = ',
41                             opt_nei)
42                         print('      Highest CH_score is: ',opt_ch, 'Corresponding
43 SH_score is: ', opt_sh)
44                         ch_score.append(opt_ch)
45                         sh_score.append(opt_sh)
46                         real_cluster.append(n_cluster)
47
48 AC_clusters=plt.figure()
49 plt.plot(num_of_cluster,ch_score)
50 plt.plot(num_of_cluster,sh_score)
51 plt.plot(num_of_cluster, real_cluster)
52 plt.grid()
53 plt.legend(["ch_score", "sh_score * 100","real_cluster_num"],bbox_to_anchor=(1.1,1))
54 plt.title("CH_score with different number of clusters")
55 plt.xlabel('Number of Clusters')
56 _ = plt.show()

# parameter matrix
init_list = np.arange(5,30,5)

```

```
54 gamma_list = np.arange(0.6,1.5,0.2)
55 aff_list = ['nearest_neighbors','precomputed','rbf']
56 neighbor_list = np.arange(5,30,5)
57
58 import warnings
59 warnings.filterwarnings('ignore')
60
61 # before pca
62 SC_grid_search(train_df,20,init_list,gamma_list,aff_list,neighbor_list)
63
64 # histogram
65 sc_model_nopca = SpectralClustering(n_clusters=4, n_init= 5, gamma =
66     0.6, affinity = 'nearest_neighbors', n_neighbors = 10,
67     random_state=42)
68 sc_predicted_nopca = sc_model_nopca.fit_predict(train_df)
69 print('ch_score = %.2f' % calinski_harabaz_score(X = train_df,labels =
70     sc_predicted_nopca))
71 print('sh-score = %.2f' % silhouette_score(X = train_df, labels =
72     sc_predicted_nopca))
73 plt.hist(sc_predicted_nopca,bins=4)
74
75 # result std
76 result = pd.DataFrame(sc_predicted_nopca)
77 result[0].value_counts().std()
78
79 # histogram
80 sc_model_nopca = SpectralClustering(n_clusters=9, n_init= 10, gamma =
81     0.6, affinity = 'nearest_neighbors', n_neighbors = 10,
82     random_state=42)
83 sc_predicted_nopca = sc_model_nopca.fit_predict(train_df)
84 print('ch_score = %.2f' % calinski_harabaz_score(X = train_df,labels =
85     sc_predicted_nopca))
86 print('sh-score = %.2f' % silhouette_score(X = train_df, labels =
87     sc_predicted_nopca))
88 plt.hist(sc_predicted_nopca, bins =9)
89
90 # result std
91 result = pd.DataFrame(sc_predicted_nopca)
92 result[0].value_counts().std()
93
94 # after std with outliers
95 SC_grid_search(m_pca,20,init_list,gamma_list,aff_list,neighbor_list)
96
97 # histogram
98 sc_model_pca = SpectralClustering(n_clusters=4, n_init= 5, gamma = 0.6,
99     affinity = 'nearest_neighbors', n_neighbors = 10, random_state
100    =42)
101 sc_predicted_pca = sc_model_pca.fit_predict(m_pca)
102 print('ch_score = %.2f' % calinski_harabaz_score(X = m_pca,labels =
103     sc_predicted_pca))
104 print('sh-score = %.2f' % silhouette_score(X = m_pca, labels =
105     sc_predicted_pca))
106 plt.hist(sc_predicted_pca,bins=4)
107
108 # result std
109 result = pd.DataFrame(sc_predicted_pca)
110 result[0].value_counts().std()
111
112 # histogram
113 sc_model_pca = SpectralClustering(n_clusters=8, n_init= 20, gamma =
114     0.6, affinity = 'nearest_neighbors', n_neighbors = 15,
115     random_state=42)
116 sc_predicted_pca = sc_model_pca.fit_predict(m_pca)
```

```

103 print('ch_score = %.2f' % calinski_harabaz_score(X = m_pca, labels =
104     sc_predicted_pca))
105 print('sh-score = %.2f' % silhouette_score(X = m_pca, labels =
106     sc_predicted_pca))
107 plt.hist(sc_predicted_pca,bins=8)
108
109 # result std
110 result = pd.DataFrame(sc_predicted_pca)
111 result[0].value_counts().std()
112
113 # histogram
114 sc_model_pca = SpectralClustering(n_clusters=12, n_init= 10, gamma =
115     0.6, affinity = 'nearest_neighbors', n_neighbors = 15,
116     random_state=42)
117 sc_predicted_pca = sc_model_pca.fit_predict(m_pca)
118 print('ch_score = %.2f' % calinski_harabaz_score(X = m_pca, labels =
119     sc_predicted_pca))
120 print('sh-score = %.2f' % silhouette_score(X = m_pca, labels =
121     sc_predicted_pca))
122 plt.hist(sc_predicted_pca,bins=12)
123
124 # result std
125 result = pd.DataFrame(sc_predicted_pca)
126 result[0].value_counts().std()
127
128 # after pca without outliers
129 SC_grid_search(pca_wo_o,20,init_list,gamma_list,aff_list,neighbor_list)
130
131 # histogram
132 sc_model_wo_o = SpectralClustering(n_clusters=8, n_init= 5, gamma =
133     0.6, affinity = 'nearest_neighbors', n_neighbors = 25,
134     random_state=42)
135 sc_predicted_wo_o = sc_model_wo_o.fit_predict(pca_wo_o)
136 print('ch_score = %.2f' % calinski_harabaz_score(X = pca_wo_o, labels =
137     sc_predicted_wo_o))
138 print('sh-score = %.2f' % silhouette_score(X = pca_wo_o, labels =
139     sc_predicted_wo_o))
140 plt.hist(sc_predicted_wo_o,bins=8)
141
142 # result std
143 result = pd.DataFrame(sc_predicted_wo_o)
144 result[0].value_counts().std()
145
146 # labelling outliers
147 pca_wo_o3 = pca_wo_o.copy()
148 pca_wo_o3['label'] = sc_predicted_wo_o
149
150 sc_center = pd.DataFrame()
151 for i in range(8):
152     sc_center[i] = np.mean(pca_wo_o3[pca_wo_o3['label']==i],axis=0)
153
154 sc_center = sc_center.drop(['label']).transpose()
155
156 sc_o_labels = classify_outliers(sc_center,pca_o)
157 print(calinski_harabaz_score(X =pd.concat([pca_wo_o,pca_o]) ,labels =
158     list(sc_predicted_wo_o)+sc_o_labels ))
159 plt.hist(list(sc_predicted_wo_o)+sc_o_labels,bins=8)

```

B.6 Model: Mean Shift

```

1  from sklearn.cluster import MeanShift
2  import sklearn.metrics
3  from sklearn.metrics import silhouette_score
4  from sklearn.metrics import davies_bouldin_score
5
6
7  MS_cluster = MeanShift(bandwidth=3.3, bin_seeding = True, cluster_all=
8      True)
9  predicted = MS_cluster.fit_predict(train_df )
10
11 from collections import Counter
12 result = Counter(predicted)
13 print(result)
14 plt.hist(predicted)
15
16 MS_cluster = MeanShift(bandwidth=3.3, bin_seeding = True, cluster_all=
17      True)
18 predicted = MS_cluster.fit_predict(data_wo_o )
19
20 from collections import Counter
21 result = Counter(predicted)
22 print(result)
23
24 # before pca with outliers
25 def find_optimal_bandwidth(start ,end ,interval):
26     num_of_cluster = []
27     ch_score = []
28     s_score = []
29     x_axis = []
30     for i in np.arange(start ,end ,interval):
31         try:
32             MS_cluster = MeanShift(bandwidth=i, bin_seeding = True ,
33             cluster_all=True)
34             predicted = MS_cluster.fit_predict(train_df )
35             ch_score.append(sklearn.metrics.calinski_harabaz_score(X =
36             train_df ,labels = predicted ))
37             s_score.append(100*silhouette_score(X = train_df ,labels =
38             predicted))
39         except:
40             continue
41
42         x_axis.append(i)
43         num_of_cluster.append(max(predicted)+1)
44 MS_cluster_bandwidth=plt.figure()
45 plt.plot(x_axis,num_of_cluster)
46 plt.plot(x_axis,ch_score)
47 plt.plot(x_axis,s_score)
48 plt.legend(["Number of clusters","CH_score","100 * Silhouette_score
49 "],bbox_to_anchor=(1.1,1))
50 plt.grid()
51 plt.title("Number of clusters silhouette_score and CH_score with
52 different bandwidth")
53 plt.xlabel('Bandwidth')
54 plt.ylabel('Number of Clusters/Score')
55 _ = plt.show()

56 print("Optimal Bandwidth:")
57 print(np.arange(start ,end ,interval)[np.argmax(ch_score)])
58 print("Optimal number for clusters:")
59 print(num_of_cluster[np.argmax(ch_score)])

```

```

54     print("Largest CH_score:")
55     print(max(ch_score))
56
57     return [ch_score, num_of_cluster]
58
59 ch_score_ms_2, num_of_cluster_ms_2 = find_optimal_bandwidth(1.5, 4.5, 0.1)
60
61 from collections import Counter
62 MS_cluster = MeanShift(bandwidth=2.86, bin_seeding = True, cluster_all=
   True)
63 predicted = MS_cluster.fit_predict(train_df )
64
65 result = Counter(predicted)
66 print(result)
67 plt.hist(predicted)
68
69 # before pca without outliers
70 def find_optimal_bandwidth(start, end, interval):
71     num_of_cluster = []
72     ch_score = []
73     s_score = []
74     x_axis = []
75     for i in np.arange(start, end, interval):
76         try:
77             MS_cluster = MeanShift(bandwidth=i, bin_seeding = True,
   cluster_all=True)
78             predicted = MS_cluster.fit_predict(data_wo_o )
79             ch_score.append(sklearn.metrics.calinski_harabaz_score(X =
   data_wo_o , labels = predicted ))
80             s_score.append(100*silhouette_score(X = data_wo_o , labels =
   predicted))
81         except:
82             continue
83
84         x_axis.append(i)
85         num_of_cluster.append(max(predicted)+1)
86 MS_cluster_bandwidth=plt.figure()
87 plt.plot(x_axis, num_of_cluster)
88 plt.plot(x_axis, ch_score)
89 plt.plot(x_axis, s_score)
90 plt.legend(["Number of clusters","CH_score","100 * Silhouette score
"], bbox_to_anchor=(1.1,1))
91 plt.grid()
92 plt.title("Number of clusters silhouette_score and CH_score with
 different bandwidth")
93 plt.xlabel('Bandwidth')
94 plt.ylabel('Number of Clusters/Score')
95 _ = plt.show()
96
97 print("Optimal Bandwidth:")
98 print(np.arange(start, end, interval)[np.argmax(ch_score)])
99 print("Optimal number for clusters:")
100 print(num_of_cluster[np.argmax(ch_score)])
101 print("Largest CH_score:")
102 print(max(ch_score))
103
104 return [ch_score, num_of_cluster]
105
106 ch_score_ms_2, num_of_cluster_ms_2 = find_optimal_bandwidth
   (1.5, 4.5, 0.01)
107
108 MS_cluster = MeanShift(bandwidth=2.75, bin_seeding = True, cluster_all=
   True)

```

```

109 predicted = MS_cluster.fit_predict(data_wo_o )
110
111 result = Counter(predicted)
112 print(result)
113 plt.hist(predicted)
114
115 # after pca with outliers
116 def find_optimal_bandwidth(start,end,interval):
117     num_of_cluster = []
118     ch_score = []
119     s_score = []
120     x_axis = []
121     for i in np.arange(start,end,interval):
122         try:
123             MS_cluster = MeanShift(bandwidth=i,bin_seeding = True,
124             cluster_all=True)
125             predicted = MS_cluster.fit_predict(pca_df )
126             ch_score.append(sklearn.metrics.calinski_harabaz_score(X =
127             pca_df ,labels = predicted ))
128             s_score.append(100*silhouette_score(X = pca_df ,labels =
129             predicted))
130         except:
131             continue
132
133         x_axis.append(i)
134         num_of_cluster.append(max(predicted)+1)
135 MS_cluster_bandwidth=plt.figure()
136 plt.plot(x_axis,num_of_cluster)
137 plt.plot(x_axis,ch_score)
138 plt.plot(x_axis,s_score)
139 plt.legend(["Number of clusters","CH_score","100 * Silhouette score
140 "],bbox_to_anchor=(1.1,1))
141 plt.grid()
142 plt.title("Number of clusters silhouette_score and CH_score with
143 different bandwidth")
144 plt.xlabel('Bandwidth')
145 plt.ylabel('Number of Clusters/Score')
146 _ = plt.show()

147 print("Optimal Bandwidth:")
148 print(np.arange(start,end,interval)[np.argmax(ch_score)])
149 print("Optimal number for clusters:")
150 print(num_of_cluster[np.argmax(ch_score)])
151 print("Largest CH_score:")
152 print(max(ch_score))

153 return [ch_score,num_of_cluster]

154 ch_score_ms_2,num_of_cluster_ms_2 = find_optimal_bandwidth(3,10,0.1)
155
156 MS_cluster = MeanShift(bandwidth=5,bin_seeding = True,cluster_all=True)
157 predicted = MS_cluster.fit_predict(pca_df )

158 result = Counter(predicted)
159 print(result)
160 plt.hist(predicted)

161 # after pca without outliers
162 def find_optimal_bandwidth(start,end,interval):
163     num_of_cluster = []
164     ch_score = []
165     s_score = []
166     x_axis = []

```

```

167     for i in np.arange(start,end,interval):
168         try:
169             MS_cluster = MeanShift(bandwidth=i,bin_seeding = True,
170             cluster_all=True)
171             predicted = MS_cluster.fit_predict(pca_wo_o )
172             ch_score.append(sklearn.metrics.calinski_harabaz_score(X =
173             pca_wo_o ,labels = predicted ))
174             s_score.append(100*silhouette_score(X = pca_wo_o ,labels =
175             predicted))
176             except:
177                 continue
178
179             x_axis.append(i)
180             num_of_cluster.append(max(predicted)+1)
181             MS_cluster_bandwidth=plt.figure()
182             plt.plot(x_axis,num_of_cluster)
183             plt.plot(x_axis,ch_score)
184             plt.plot(x_axis,s_score)
185             plt.legend(["Number of clusters","CH_score","100 * Silhouette_score
186             "],bbox_to_anchor=(1.1,1))
187             plt.grid()
188             plt.title("Number of clusters silhouette_score and CH_score with
189             different bandwidth")
190             plt.xlabel('Bandwidth')
191             plt.ylabel('Number of Clusters/Score')
192             _ = plt.show()
193
194             print("Optimal Bandwidth:")
195             print(np.arange(start,end,interval)[np.argmax(ch_score)])
196             print("Optimal number for clusters:")
197             print(num_of_cluster[np.argmax(ch_score)])
198             print("Largest CH_score:")
199             print(max(ch_score))
200
201             return [ch_score,num_of_cluster]
202
203
204     ch_score_ms_2,num_of_cluster_ms_2 = find_optimal_bandwidth(3,12,0.1)

```

B.7 Model: Affinity Propagation

```

1 from sklearn.cluster import AffinityPropagation
2
3 # before pca with outliers
4 def find_optimal_damping(start,end,interval):
5     num_of_cluster = []
6     ch_score = []
7     s_score = []
8     x = []
9     for i in np.arange(start,end,interval):
10         AP_cluster = AffinityPropagation(preference=i,convergence_iter
= 50,max_iter=400)
11         predicted = AP_cluster.fit_predict(pca_df )
12         try:
13             ch_score.append(sklearn.metrics.calinski_harabaz_score(X =
pca_df ,labels = predicted ))
14         except:
15             continue
16         x.append(i)
17         num_of_cluster.append(len(set(predicted) ))
18
19         s_score.append(100*silhouette_score(X = pca_df ,labels =
predicted))
20     AP_cluster_damping=plt.figure()
21     plt.plot(x,num_of_cluster)
22     plt.plot(x,ch_score)
23     plt.plot(x,s_score)
24
25     plt.legend(["Number of clusters","CH_score","Silhouette_score"],bbox_to_anchor=(1.1,1))
26     plt.grid()
27     plt.title("Number of clusters,Silhouette_score and CH_score with
different performance")
28     plt.xlabel('damping')
29     plt.ylabel('Number of Clusters')
30     _ = plt.show()
31
32     print("Optimal preference:")
33     print(x[np.argmax(ch_score)])
34     print("Optimal number for clusters:")
35     print(num_of_cluster[np.argmax(ch_score)])
36     print("Largest CH_score:")
37     print(max(ch_score))
38
39     return [ch_score,num_of_cluster]
40
41 ch_score_ms,num_of_cluster_ms = find_optimal_damping(-600,-60,10)
42
43 for item in [-500,-450,-400,-200]:
44     MS_cluster = AffinityPropagation(preference = item,convergence_iter
= 50,max_iter=400)
45     predicted =MS_cluster.fit_predict(train_df)
46     print("Number of Clusters:{}".format(max(predicted)+1))
47
48     from collections import Counter
49     result = Counter(predicted)
50     plt.figure()
51     plt.hist(predicted,bins= range(max(predicted)+1))
52     plt.show()
53
54 # before pca without outliers

```

```

55 def find_optimal_damping(start,end,interval):
56     num_of_cluster = []
57     ch_score = []
58     s_score = []
59     x = []
60     for i in np.arange(start,end,interval):
61         AP_cluster = AffinityPropagation(preference=i,convergence_iter
62 = 50,max_iter=400)
63         predicted = AP_cluster.fit_predict(data_wo_o )
64         try:
65             ch_score.append(sklearn.metrics.calinski_harabaz_score(X =
66 data_wo_o ,labels = predicted ))
67         except:
68             continue
69         x.append(i)
70         num_of_cluster.append(len(set(predicted) ))
71
72         s_score.append(100*silhouette_score(X = data_wo_o ,labels =
73 predicted))
74         AP_cluster_damping=plt.figure()
75         plt.plot(x,num_of_cluster)
76         plt.plot(x,ch_score)
77         plt.plot(x,s_score)
78
79         plt.legend(["Number of clusters","CH_score","Silhouette_score"],bbox_to_anchor=(1.1,1))
80         plt.grid()
81         plt.title("Number of clusters,Silhouette_score and CH_score with
82 different performance")
83         plt.xlabel('damping')
84         plt.ylabel('Number of Clusters')
85         _ = plt.show()
86
87         print("Optimal performance:")
88         print(x[np.argmax(ch_score)])
89         print("Optimal number for clusters:")
90         print(num_of_cluster[np.argmax(ch_score)])
91         print("Largest CH_score:")
92         print(max(ch_score))
93
94         return [ch_score,num_of_cluster]
95
96 ch_score_ms,num_of_cluster_ms = find_optimal_damping(-600,-60,10)
97
98 for item in [-300,-240,-150,-100]:
99     MS_cluster = AffinityPropagation(preference = item,convergence_iter
100 = 50,max_iter=400)
101     predicted =MS_cluster.fit_predict(data_wo_o )
102     print("Number of Clusters:{}".format(max(predicted)+1))
103
104     from collections import Counter
105     result = Counter(predicted)
106     plt.figure()
107     plt.hist(predicted,bins= range(max(predicted)+1))
108     plt.show()
109
110     AP_cluster = AffinityPropagation(preference = -150,convergence_iter =
111         50,max_iter=400)
112     predicted = AP_cluster.fit_predict(data_wo_o )
113     print("Number of Clusters:{}".format(max(predicted)))
114
115     from collections import Counter
116     result = Counter(predicted)

```

```

111 plt.figure()
112 plt.hist(predicted,bins= range(max(predicted)+1))
113 plt.show()
114 display(sorted(result.items(), key=lambda x: x[0]))
115 print("CH score = {}".format(sklearn.metrics.calinski_harabaz_score(X =
116     data_wo_o ,labels = predicted )))
117
118 # outliers
119 from scipy.spatial import distance
120 #center can be list or array
121 #assume each center correponds to label of [0,1,2...]
122 #return a list of label
123 def classify_outliers(centers,outliers_df):
124     labels = []
125     for sample in np.array(outliers_df):
126         min_center = float('inf')
127         min_center_th = -10086
128         for i,center in enumerate(centers):
129             dst = distance.euclidean(center,sample)
130             if dst<min_center:
131                 min_center = dst
132                 min_center_th = i
133         labels.append(min_center_th)
134     return labels
135
136 out_labels = classify_outliers(AP_cluster.cluster_centers_ ,data_o)
137
138 sklearn.metrics.calinski_harabaz_score(X =pd.concat([data_wo_o ,data_o])
139     ,labels = list(predicted)+out_labels )
140
141 # after pca with outliers
142 def find_optimal_damping(start,end,interval):
143     num_of_cluster = []
144     ch_score = []
145     s_score = []
146     x = []
147     for i in np.arange(start,end,interval):
148         AP_cluster = AffinityPropagation(preference=i,convergence_iter
149 = 50,max_iter=400)
150         predicted = AP_cluster.fit_predict(pca_df )
151         try:
152             ch_score.append(sklearn.metrics.calinski_harabaz_score(X =
153             pca_df ,labels = predicted ))
154             except:
155                 continue
156             x.append(i)
157             num_of_cluster.append(len(set(predicted) ))
158
159             s_score.append(100*silhouette_score(X = pca_df ,labels =
160             predicted))
161             AP_cluster_damping=plt.figure()
162             plt.plot(x,num_of_cluster)
163             plt.plot(x,ch_score)
164             plt.plot(x,s_score)
165
166             plt.legend(["Number of clusters","CH_score","Silhouette_score"],bbox_to_anchor=(1.1,1))
167             plt.grid()
168             plt.title("Number of clusters,Silhouette_score and CH_score with
169             different performance")
170             plt.xlabel('damping')
171             plt.ylabel('Number of Clusters')

```

```

167     _ = plt.show()
168
169     print("Optimal performance:")
170     print(x[np.argmax(ch_score)])
171     print("Optimal number for clusters:")
172     print(num_of_cluster[np.argmax(ch_score)])
173     print("Largest CH_score:")
174     print(max(ch_score))
175
176     return [ch_score, num_of_cluster]
177
178 ch_score_ms, num_of_cluster_ms = find_optimal_damping(-400, -60, 10)
179
180 AP_cluster = AffinityPropagation(preference = -300, convergence_iter =
181                                     50, max_iter=400)
182 predicted = AP_cluster.fit_predict(pca_df)
183 print("Number of Clusters:{}".format(max(predicted)))
184
185 from collections import Counter
186 result = Counter(predicted)
187
188 plt.figure()
189 plt.hist(predicted, bins= max(predicted)+1)
190 plt.show()
191 display(sorted(result.items(), key=lambda x: x[0]))
192
193 # after pca without outliers
194 def find_optimal_damping(start, end, interval):
195     num_of_cluster = []
196     ch_score = []
197     s_score = []
198     x = []
199     for i in np.arange(start, end, interval):
200         AP_cluster = AffinityPropagation(preference=i, convergence_iter =
201                                         50, max_iter=400)
202         predicted = AP_cluster.fit_predict(pca_wo_o )
203         try:
204             ch_score.append(sklearn.metrics.calinski_harabaz_score(X =
205                         pca_wo_o , labels = predicted ))
206             except:
207                 continue
208             x.append(i)
209             num_of_cluster.append(len(set(predicted) ))
210
211             s_score.append(100*silhouette_score(X = pca_wo_o , labels =
212                                         predicted))
213             AP_cluster_damping=plt.figure()
214             plt.plot(x,num_of_cluster)
215             plt.plot(x,ch_score)
216             plt.plot(x,s_score)
217
218             plt.legend(["Number of clusters","CH_score","Silhouette_score"],
219             bbox_to_anchor=(1.1,1))
220             plt.grid()
221             plt.title("Number of clusters, Silhouette_score and CH_score with
222             different performance")
223             plt.xlabel('damping')
224             plt.ylabel('Number of Clusters')
225             _ = plt.show()
226
227             print("Optimal performance:")
228             print(x[np.argmax(ch_score)])
229             print("Optimal number for clusters:")

```

```

224 print(num_of_cluster[np.argmax(ch_score)])
225 print("Largest CH_score:")
226 print(max(ch_score))
227
228 return [ch_score, num_of_cluster]
229
230 ch_score_ms, num_of_cluster_ms = find_optimal_damping(-400, -60, 10)
231
232 AP_cluster = AffinityPropagation(preference = -125, convergence_iter =
233 50, max_iter=400)
234 predicted = AP_cluster.fit_predict(pca_wo_o)
235 print("Number of Clusters:{}".format(max(predicted)))
236
237 from collections import Counter
238 result = Counter(predicted)
239
240 plt.figure()
241 plt.hist(predicted, bins= max(predicted)+1)
242 plt.show()
243 display(sorted(result.items(), key=lambda x: x[0]))
244 print("CH score = {}".format(sklearn.metrics.calinski_harabaz_score(X =
245 pca_wo_o ,labels = predicted)))
246 print("SH score = {}".format(silhouette_score(X = pca_wo_o ,labels =
247 predicted)))
248
249 # labelling outliers
250 from scipy.spatial import distance
251 #center can be list or array
252 #assume each center correponds to label of [0,1,2...]
253 #return a list of label
254 def classify_outliers(centers,outliers_df):
255     labels = []
256     for sample in np.array(outliers_df):
257         min_center = float('inf')
258         min_center_th = -10086
259         for i,center in enumerate(centers):
260             dst = distance.euclidean(center,sample)
261             if dst<min_center:
262                 min_center = dst
263                 min_center_th = i
264         labels.append(min_center_th)
265     return labels
266
267 out_labels = classify_outliers(AP_cluster.cluster_centers_ ,pca_o)
268 sklearn.metrics.calinski_harabaz_score(X =pd.concat([pca_wo_o,pca_o])
269 ,labels = list(predicted)+out_labels )

```

B.8 Model: DBSCAN

```

1  from sklearn.cluster import DBSCAN
2
3  # use data after pca
4  for eps in np.arange(1,25,0.5):
5      MS_cluster = DBSCAN(eps=eps,metric ="l1",min_samples = 1)
6      predicted =MS_cluster.fit_predict(data_wo_o)
7      print("Number of Clusters:{}".format(max(predicted)))
8      from collections import Counter
9      result = Counter(predicted)
10     plt.figure
11     plt.hist(predicted)
12     plt.show()
13
14
15 def find_optimal_eps(start,end,interval,min_sample,metric):
16     num_of_cluster = []
17     ch_score = []
18     s_score = []
19     x_axis = []
20     for i in np.arange(start,end,interval):
21
22         DBSCAN_cluster = DBSCAN(eps=i,metric = metric,min_samples =
23 min_sample)
24         predicted = DBSCAN_cluster.fit_predict(pca_wo_o )
25         try:
26             ch_score.append(sklearn.metrics.calinski_harabaz_score(X =
27 pca_wo_o ,labels = predicted ))
28         except:
29             continue
30         s_score.append(10*silhouette_score(X = pca_wo_o ,labels =
31 predicted))
32         x_axis.append(i)
33         num_of_cluster.append(len(set(predicted) ))
34
35
36 DBSCAN_cluster_eps=plt.figure()
37
38 plt.plot(x_axis,num_of_cluster)
39 plt.plot(x_axis,ch_score)
40 plt.plot(x_axis,s_score)
41 plt.legend(["Number of clusters","CH_score","Silhouette_score"],bbox_to_anchor=(1.1,1))
42 plt.title("Number of clusters,Silhouette_score and CH_score with
43 different damping")
44 plt.grid()
45 plt.title("Number of clusters and CH_score with different eps")
46 plt.xlabel('eps')
47 plt.ylabel('Number of Clusters/CH_score')
48 _ = plt.show()
49
50
51 print("Optimal eps:")
52 print(np.arange(start,end,interval)[np.argmax(ch_score)])
53 print("Optimal number for clusters:")
54 print(num_of_cluster[np.argmax(ch_score)])
55 print("Largest CH_score:")
56 print(max(ch_score))
57
58 return [ch_score,num_of_cluster]
59

```

```

55 for metric in ['cityblock', 'cosine', 'euclidean', 'l1', 'l2', '',
56     manhattan']:
57     print("")
58     print("Distance Metric = {}".format(metric))
59     for min_sample in np.arange(1,15,1):
60         print("")
61         print("min_sample = {}".format(min_sample))
62         try:
63             ch_score_DBSCAN, num_of_cluster_DBSCAN = find_optimal_eps
64             (1,8,0.01,min_sample,metric)
65         except:
66             break
67
68 def find_optimal_min_samples(start,end,interval):
69     num_of_cluster = []
70     ch_score = []
71     s_score = []
72     for i in np.arange(start,end,interval):
73         DBSCAN_cluster = DBSCAN(eps = 4,min_samples=i,metric =
74         'manhattan')
75         predicted = DBSCAN_cluster.fit_predict(pca_wo_o )
76         num_of_cluster.append(len(set(predicted) ))
77         ch_score.append(sklearn.metrics.calinski_harabaz_score(X =
78         pca_wo_o ,labels = predicted ))
79         s_score.append(10*silhouette_score(X = pca_wo_o ,labels =
80         predicted))
81         DBSCAN_cluster_min_samples=plt.figure()
82         plt.plot(np.arange(start,end,interval),num_of_cluster)
83         plt.plot(np.arange(start,end,interval),ch_score)
84         plt.plot(np.arange(start,end,interval),s_score)
85         plt.legend(["Number of clusters","CH_score","Silhouette_score"],
86         bbox_to_anchor=(1.1,1))
87         plt.grid()
88         plt.title("Number of clusters, Silhouette_score and CH_score with
89         different min_samples")
90         plt.xlabel('min_samples')
91         plt.ylabel('Number of Clusters/CH_score')
92         _ = plt.show()
93
94     print("Optimal min_samples:")
95     print(np.arange(start,end,interval)[np.argmax(ch_score)])
96     print("Optimal number for clusters:")
97     print(num_of_cluster[np.argmax(ch_score)])
98     print("Largest CH_score:")
99     print(max(ch_score))
100
101     return [ch_score,num_of_cluster]
102
103 _ = find_optimal_min_samples(2,13,1)
104
105 def find_optimal_grid():
106     num_of_cluster = []
107     ch_score = []
108     for i in np.arange(1,300,1):
109         cluster_temp = []
110         ch_score_temp = []
111         for j in np.arange(1,100,0.1):
112             DBSCAN_cluster = DBSCAN(eps = j,min_samples=i,metric =
113             "cosine")
114             try:
115                 predicted = DBSCAN_cluster.fit_predict(pca_wo_o )
116                 cluster_temp.append(max(predicted)+1 )
117             except:
118                 break
119         num_of_cluster.append(len(set(cluster_temp)))
120         ch_score.append(calinski_harabaz_score(X = pca_wo_o ,
121         cluster_temp))
122
123     return [num_of_cluster,ch_score]
124
125 _ = find_optimal_grid()

```

```

109         ch_score_temp.append(sklearn.metrics.
110     calinski_harabaz_score(X = pca_wo_o ,labels = predicted ))
111     except:
112         cluster_temp.append(0)
113         ch_score_temp.append(0)
114
115     num_of_cluster.append(cluster_temp)
116     ch_score.append(ch_score_temp)
117
118 ch_score = np.array(ch_score)
119 num_of_cluster = np.array(num_of_cluster)
120
121 k = np.argmax(ch_score)
122 count = 0
123
124 for i in np.arange(2,30,1):
125     for j in np.arange(1,10,0.1):
126         if count==k:
127             eps_o = j
128             min_samples_o = i
129         count+=1
130
131
132
133 print("Optimal min_samples:")
134 print(min_samples_o)
135 print("Optimal eps:")
136 print(eps_o)
137
138 DBSCAN_cluster = DBSCAN(eps = eps_o,min_samples=min_samples_o)
139 predicted = DBSCAN_cluster.fit_predict(pca_wo_o )
140 clusters_o = max(predicted)+1
141 print("Optimal number for clusters:")
142 print(clusters_o)
143 print("Largest CH_score:")
144 print(np.max(ch_score))
145
146 return [ch_score,num_of_cluster]
147
148 _ = find_optimal_grid()
149
150 _[0]
```

Bibliography

- [1] G. Dong and H. Liu, *Feature Engineering for Machine Learning and Data Analytics*. Boca Raton, FL, USA: CRC Press, Inc., 1st ed., 2018.
- [2] “Stock prediction with ml: Feature engineering.” http://alphascientist.com/feature_engineering.html. Accessed: 2018-07-30.
- [3] G. S. Bonde, “Extracting the best features from multi-company stock data to improve stock price prediction,” 2012.
- [4] Wikipedia contributors, “Principal component analysis — Wikipedia, the free encyclopedia,” 2019. [Online; accessed 12-August-2019].
- [5] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Pas-sos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, “Scikit-learn: Machine learning in Python,” *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [6] A. Singh, A. Yadav, and A. Rana, “K-means with three different distance metrics,” *International Journal of Computer Applications*, vol. 67, no. 10, 2013.
- [7] M. Charrad, N. Ghazzali, V. Boiteau, and A. Niknafs, “Nbclust package: finding the relevant number of clusters in a dataset,” *User!* 2012, 2012.
- [8] L. Kaufman and P. J. Rousseeuw, *Finding groups in data: an introduction to cluster analysis*, vol. 344. John Wiley & Sons, 2009.
- [9] M. S. Babu, N. Geethanjali, and B. Satyanarayana, “Clustering approach to stock market prediction,” *International Journal of Advanced Networking and Applications*, vol. 3, no. 4, p. 1281, 2012.
- [10] T. Wittman, “Time-series clustering and association analysis of financial data,” *University of Texas, Austin*, 2002.
- [11] M. Gavrilov, D. Anguelov, P. Indyk, and R. Motwani, “Mining the stock market: Which measure is best,” in *Proc. of the 6th ACM SIGKDD*, pp. 487–496, 2000.
- [12] F. Cai, N.-A. Le-Khac, and T. Kechadi, “Clustering approaches for financial data analysis: a survey,” *arXiv preprint arXiv:1609.08520*, 2016.
- [13] M. Momeni, M. Mohseni, and M. Soofi, “Clustering stock market companies via k-means algorithm,” *Kuwait Chapter of Arabian Journal of Business and Management Review*, vol. 33, no. 2578, pp. 1–10, 2015.
- [14] ttnphns (<https://stats.stackexchange.com/users/3277/ttnphns>), “What is an acceptable value of the calinski harabasz (ch) criterion?.” Cross Validated. URL:<https://stats.stackexchange.com/q/52863> (version: 2019-07-15).

- [15] K. Sasirekha and P. Baby, "Agglomerative hierarchical clustering algorithm-a," *International Journal of Scientific and Research Publications*, vol. 83, p. 83, 2013.